

Statistics and Information Theory

John Duchi

October 22, 2024

Contents

1	Introduction and setting	10
1.1	Information theory	10
1.2	Moving to statistics and machine learning	11
1.3	Outline and chapter discussion	12
1.4	A remark about measure theory	14
2	An information theory review	15
2.1	Basics of Information Theory	15
2.1.1	Definitions	15
2.1.2	Chain rules and related properties	20
2.1.3	Data processing inequalities:	22
2.2	General divergence measures and definitions	23
2.2.1	Partitions, algebras, and quantizers	23
2.2.2	KL-divergence	24
2.2.3	f -divergences	26
2.2.4	Inequalities and relationships between divergences	28
2.2.5	Convexity and data processing for divergence measures	32
2.3	First steps into optimal procedures: testing inequalities	33
2.3.1	Le Cam's inequality and binary hypothesis testing	33
2.3.2	Fano's inequality and multiple hypothesis testing	35
2.4	A first operational result: entropy and source coding	37
2.4.1	The source coding problem	37
2.4.2	The Kraft-McMillan inequalities	38
2.4.3	Entropy rates and longer codes	41
2.5	Bibliography	43
2.6	Exercises	43
3	Exponential families and statistical modeling	48
3.1	Exponential family models	48
3.2	Why exponential families?	50
3.2.1	Fitting an exponential family model	53
3.3	Divergence measures and information for exponential families	54
3.4	Generalized linear models and regression	55
3.4.1	Fitting a generalized linear model from a sample	58
3.4.2	The information in a generalized linear model	59
3.5	Lower bounds on testing a parameter's value	61

3.6	Deferred proofs	63
3.6.1	Proof of Proposition 3.2.2	63
3.7	Bibliography	65
3.8	Exercises	65
I	Concentration, information, stability, and generalization	66
4	Concentration Inequalities	67
4.1	Basic tail inequalities	67
4.1.1	Sub-Gaussian random variables	69
4.1.2	Sub-exponential random variables	73
4.1.3	Orlicz norms	77
4.1.4	First applications of concentration: random projections	80
4.1.5	A second application of concentration: codebook generation	81
4.2	Martingale methods	83
4.2.1	Sub-Gaussian martingales and Azuma-Hoeffding inequalities	84
4.2.2	Examples and bounded differences	85
4.3	Matrix concentration	88
4.4	Technical proofs	91
4.4.1	Proof of Theorem 4.1.11	91
4.4.2	Proof of Theorem 4.1.15	92
4.4.3	Proof of Theorem 5.1.6	93
4.4.4	Proof of Proposition 4.3.2	93
4.5	Bibliography	95
4.6	Exercises	95
5	Estimation and generalization	102
5.1	Uniformity and metric entropy	102
5.1.1	Symmetrization and uniform laws	102
5.1.2	Metric entropy, coverings, and packings	106
5.1.3	Application: matrix concentration	109
5.2	Generalization bounds	110
5.2.1	Finite and countable classes of functions	112
5.2.2	Large classes	113
5.2.3	Structural risk minimization and adaptivity	116
5.3	M-estimators and estimation	118
5.3.1	Standard conditions and convex optimization	119
5.3.2	Some growth properties of convex functions	120
5.3.3	Convergence analysis for convex M-estimators	122
5.3.4	Consequences for exponential families and generalized linear models	124
5.3.5	Proof of Theorem 5.3.8	126
5.4	Exercises	127

6	Generalization and stability	132
6.1	The variational representation of Kullback-Leibler divergence	133
6.2	PAC-Bayes bounds	134
6.2.1	Relative bounds	137
6.2.2	A large-margin guarantee	139
6.2.3	A mutual information bound	141
6.3	Interactive data analysis	142
6.3.1	The interactive setting	143
6.3.2	Second moment errors and mutual information	144
6.3.3	Limiting interaction in interactive analyses	145
6.3.4	Error bounds for a simple noise addition scheme	150
6.4	Bibliography and further reading	152
6.5	Exercises	152
7	Advanced concentration inequalities	157
7.1	From divergences to concentration and back	157
7.1.1	Concentration of covariance matrices via the variational representation	159
7.1.2	A generalized connection between moment generating functions and divergence	161
7.2	Transportation inequalities	163
7.2.1	A tensorized transportation inequality	165
7.2.2	A heuristic proof of Theorem 7.2.1	166
7.2.3	Proof of Corollary 7.2.2	167
7.3	Some applications of concentration and the variational inequality	168
7.3.1	Metric Gaussianity, transport inequalities, and expansion of sets	168
7.3.2	A weak and strong converse for hypothesis testing	171
7.4	Discussion and bibliographic remarks	174
7.5	Exercises	174
8	Privacy and disclosure limitation	177
8.1	Disclosure limitation, privacy, and definitions	177
8.1.1	Basic mechanisms	179
8.1.2	Resilience to side information, Bayesian perspectives, and data processing . .	183
8.2	Weakenings of differential privacy	185
8.2.1	Basic mechanisms	186
8.2.2	Connections between privacy measures	188
8.2.3	Side information protections under weakened notions of privacy	191
8.3	Composition and privacy based on divergence	194
8.3.1	Composition of Rényi-private channels	194
8.3.2	Privacy games and composition	195
8.4	Additional mechanisms and privacy-preserving algorithms	197
8.4.1	The exponential mechanism	197
8.4.2	Local sensitivities and the inverse sensitivity mechanism	200
8.5	Deferred proofs	205
8.5.1	Proof of Lemma 8.2.10	205
8.6	Bibliography	208
8.7	Exercises	208

II	Fundamental limits and optimality	215
9	Minimax lower bounds: the Le Cam, Fano, and Assouad methods	217
9.1	Basic framework and minimax risk	217
9.2	Preliminaries on methods for lower bounds	219
9.2.1	From estimation to testing	220
9.2.2	Inequalities between divergences and product distributions	221
9.2.3	Metric entropy and packing numbers	223
9.3	Le Cam's method	224
9.4	Fano's method	226
9.4.1	The classical (local) Fano method	226
9.4.2	A distance-based Fano method	231
9.5	Assouad's method	235
9.5.1	Well-separated problems	235
9.5.2	From estimation to multiple binary tests	235
9.5.3	Example applications of Assouad's method	237
9.6	Deferred proofs	239
9.6.1	Proof of Proposition 9.4.6	239
9.6.2	Proof of Corollary 9.4.7	240
9.6.3	Proof of Lemma 9.5.2	240
9.7	Bibliography	241
9.8	Exercises	241
10	Beyond local minimax techniques	250
10.1	Nonparametric regression: minimax upper and lower bounds	250
10.1.1	Kernel estimates of the function	251
10.1.2	Minimax lower bounds on estimation with Assouad's method	253
10.2	Global Fano Method	256
10.2.1	A mutual information bound based on metric entropy	256
10.2.2	Minimax bounds using global packings	258
10.2.3	Example: non-parametric regression	259
10.3	Strong converses and high-probability lower bounds	260
10.3.1	Refined Fano inequalities	262
10.3.2	High probability estimation lower bounds	265
10.3.3	Proof of Theorem 10.3.5	266
10.4	Exercises	268
11	Constrained risk inequalities	269
11.1	Strong data processing inequalities	269
11.2	Local privacy	272
11.3	Communication complexity	276
11.3.1	Classical communication complexity problems	276
11.3.2	Deterministic communication: lower bounds and structure	279
11.3.3	Randomization, information complexity, and direct sums	281
11.3.4	The structure of randomized communication and communication complexity of primitives	285
11.4	Communication complexity in estimation	288

11.4.1	Direct sum communication bounds	289
11.4.2	Communication data processing	290
11.4.3	Applications: communication and privacy lower bounds	292
11.5	Proof of Theorem 11.4.4	296
11.5.1	Proof of Lemma 11.5.3	300
11.6	Bibliography	301
11.7	Exercises	302
12	Squared error and asymptotically exact optimality guarantees	307
12.1	The Cramér-Rao inequality	308
12.1.1	Compact sets and the failure of the Cramér-Rao bound	309
12.1.2	Regularization and the failure of the Cramér-Rao bound	309
12.2	The van Trees inequality: a Bayesian Cramér-Rao bound	310
12.2.1	The van Trees inequality in one dimension	311
12.2.2	The van Trees inequality in d -dimensions	312
12.2.3	The van Trees inequality for a function of the parameter	314
12.3	Beyond parametric problems	318
12.3.1	An extended example: M-estimation lower bounds	321
12.4	Super-efficiency and instance optimality	323
12.5	Applications in privacy	325
12.6	Bibliography and further reading	325
12.7	Exercises	325
13	Testing and functional estimation	328
13.1	Geometrizing rates of convergence	328
13.1.1	Fisher information and divergence measures	332
13.1.2	Valid asymptotic information expansions of divergences	334
13.2	Le Cam's convex hull method	336
13.2.1	The χ^2 -mixture bound	337
13.2.2	Estimating the norm of a Gaussian vector	340
13.2.3	Lower bounds on estimating integral functionals	342
13.3	Minimax hypothesis testing	345
13.3.1	Detecting a difference in populations	346
13.3.2	Signal detection and testing a Gaussian mean	348
13.3.3	Goodness of fit and two-sample tests for multinomials	350
13.3.4	Detecting sparse signals and phase transitions	353
13.4	Instance-optimal lower bounds and super-efficiency	358
13.4.1	Risk transfer inequalities	358
13.4.2	A general risk transfer bound	362
13.4.3	Risk transfer with mixtures	363
13.5	Deferred and technical proofs	366
13.5.1	Proof of Lemma 13.1.6	366
13.5.2	Proof of Lemma 13.1.10	366
13.6	Bibliography	368
13.7	A useful divergence calculation	369
13.8	Exercises	370

III	Entropy, predictions, divergences, and information	377
14	Predictions, loss functions, and entropies	379
14.1	Proper losses, scoring rules, and generalized entropies	380
14.1.1	A convexity primer	381
14.1.2	From a proper loss to an entropy	383
14.1.3	The information in an experiment	385
14.2	Characterizing proper losses and Bregman divergences	386
14.2.1	Characterizing proper losses for Y taking finitely many vales	386
14.2.2	General proper losses	389
14.2.3	Proper losses and vector-valued Y	393
14.3	From entropies to convex losses, arbitrary predictions, and link functions	396
14.3.1	Convex conjugate linkages	396
14.3.2	Convex conjugate linkages with affine constraints	400
14.4	Exponential families, maximum entropy, and log loss	403
14.4.1	Maximizing entropy	405
14.5	Technical and deferred proofs	409
14.5.1	Finalizing the proof of Theorem 14.2.15	409
14.5.2	Proof of Proposition 14.4.1	410
14.5.3	Proof of Proposition 14.4.3	411
14.6	Exercises	412
15	Calibration and Proper Losses	416
15.1	Proper losses and calibration error	417
15.2	Measuring calibration	420
15.2.1	The impossibility of measuring calibration	420
15.2.2	Alternative calibration measures	423
15.3	Auditing and improving calibration at the population level	426
15.3.1	The post-processing gap and calibration audits for squared error	426
15.3.2	Calibration audits for losses based on conjugate linkages	428
15.3.3	A population-level algorithm for calibration	430
15.4	Calibeating: improving squared error by calibration	431
15.4.1	Proof of Theorem 15.4.1	434
15.5	Continuous and equivalent calibration measures	437
15.5.1	Calibration measures	438
15.5.2	Equivalent calibration measures	440
15.6	Deferred technical proofs	446
15.6.1	Proof of Lemma 15.2.1	446
15.6.2	Proof of Proposition 15.5.2	447
15.6.3	Proof of Lemma 15.5.4	448
15.6.4	Proof of Theorem 15.5.6	449
15.7	Bibliography	451
15.8	Exercises	452

16 Classification, Divergences, and Surrogate Risk	454
16.1 Surrogate risk consistency in binary classification	455
16.1.1 A general classification calibration result	458
16.1.2 Convex losses for binary classification	459
16.1.3 Proof of Theorem 16.1.5	460
16.2 General surrogate risk consistency	463
16.2.1 Uniform calibration	464
16.2.2 Pointwise calibration	465
16.2.3 Examples: multiclass surrogate risk consistency	466
16.3 Generalized entropies and surrogate risk consistency	468
16.3.1 Proof of Theorem 16.3.2	470
16.4 Structured prediction and generalized entropies	471
16.4.1 The failure of naive margin- and hinge-type losses	474
16.4.2 Structured prediction losses via the generalized entropy	476
16.4.3 Proof of Theorem 16.4.9	479
16.5 Universal loss equivalence and entropies	480
16.5.1 Proof of Theorem 16.5.1	483
16.5.2 Proof of Lemma 16.5.5	485
16.5.3 Proof of Lemma 16.5.7	486
16.6 Bibliography	486
16.7 Exercises	486
 IV Online game playing and compression	 490
17 Stochastic and online convex optimization	491
17.1 Preliminaries on convex optimization	492
17.2 Online convex optimization methods	493
17.2.1 Projected subgradient methods	494
17.2.2 Mirror descent-type methods	496
17.2.3 Convergence analysis of mirror descent	498
17.2.4 Instantiations of the regret guarantee	499
17.2.5 Proof of Theorem 17.2.9	501
17.3 Optimality guarantees and fundamental limits	503
17.3.1 From optimization to testing	504
17.3.2 Constructing hard classes of optimization problems	506
17.3.3 Instantiations and optimality	509
17.3.4 A lower bound for high-dimensional stochastic optimization	512
17.4 Online to batch conversions	513
17.5 More refined convergence guarantees	514
17.5.1 Proof of Proposition 17.5.1	515
17.6 Exercises	517
 18 Exploration, exploitation, and bandit problems	 523
18.1 The multi-armed bandit problem	523
18.2 Confidence-based algorithms	525
18.3 General losses and information-based bounds	529

18.3.1	An information-based regret bound	530
18.3.2	Posterior (Thompson) sampling	533
18.3.3	Information-based exploration	536
18.3.4	An extended example: linear bandits	539
18.4	Online gradient descent approaches	541
18.4.1	Some empirical comparisons	544
18.5	Minimax lower bounds	545
18.5.1	Action separation and a modulus of continuity	546
18.5.2	Assoaud's method for lower bounds	549
18.5.3	Proof of Theorem 18.5.1	552
18.6	Technical proofs	553
18.6.1	Proof of Lemma 18.2.1	553
18.7	Further notes and references	554
18.8	Exercises	556
19	Minimax games and Bayesian estimation	558
19.1	Robust Bayesian procedures and maximum entropy	559
19.1.1	A digression on min-max games	560
19.1.2	Saddle points for maximum entropy	561
19.1.3	Exponential family models as robust Bayesian procedures	561
19.2	The coding game and sequential prediction	563
19.3	Expected regret, information capacity, and redundancy	565
19.3.1	Information capacity and regret duality	566
19.3.2	Instantiations and corollaries of regret/capacity duality	569
19.3.3	Maximum generalized entropy and Robust Bayesian procedures	570
19.3.4	Proof of Lemma 19.3.1	572
19.4	Minimax strategies for regret	573
19.5	Mixture (Bayesian) strategies and redundancy	575
19.5.1	Bayesian redundancy and objective, reference, and Jeffreys priors	578
19.5.2	Heuristic calculations: normality and Theorem 19.5.1	580
19.6	Regret and capacity dualities	581
19.6.1	Duality when the domain is finite	581
19.6.2	Proof of Corollary 19.3.4	582
19.6.3	Regret/capacity duality for arbitrary domains	584
19.6.4	A formal statement of regret/capacity duality	588
19.7	Bibliographic details	589
19.8	Exercises	589
V	Appendices	592
A	Miscellaneous mathematical results	593
A.1	The roots of a polynomial	593
A.2	Measure-theoretic development of divergence measures	593
A.3	Integral convergence and completeness of probability spaces	593
A.4	Probabilistic convergence	594
A.4.1	Classical results on convergence in distribution	594

A.4.2	Assorted convergence results for probability distributions	595
A.5	Stirling approximations and entropy	598
B	Convex Analysis	600
B.1	Convex sets	600
B.1.1	Operations preserving convexity	602
B.1.2	Representation and separation of convex sets	604
B.2	Sublinear and support functions	608
B.3	Convex functions	611
B.3.1	Equivalent definitions of convex functions	612
B.3.2	Continuity properties of convex functions	614
B.3.3	Operations preserving convexity	620
B.3.4	Smoothness properties, first-order developments for convex functions, and subdifferentiability	623
B.3.5	Calculus rules of subgradients	628
C	Optimality, stability, and duality	631
C.1	Optimality conditions and stability properties	632
C.1.1	Subgradient characterizations for optimality	632
C.1.2	Stability properties of minimizers	634
C.2	Conjugacy and duality properties	639
C.2.1	Gradient dualities and the Fenchel-Young inequality	640
C.2.2	Smoothness and strict convexity of conjugates	641
C.2.3	Smooth convex functions	643
C.3	Limits at infinity of convex functions and sets	645
C.3.1	Boundedness and closedness of convex sets	646
C.3.2	Asymptotic growth and existence of minimizers	649
C.4	Saddle point theorems and min-max duality	651
C.4.1	Saddle points and convex conjugates	652
C.4.2	Min-max duality and the existence of saddle points	654
C.5	Exercises	655

Chapter 1

Introduction and setting

This book explores some of the (many) connections relating information theory, statistics, computation, and learning. Signal processing, machine learning, and statistics all revolve around extracting useful information from signals and data. In signal processing and information theory, a central question is how to best *design* signals—and the channels over which they are transmitted—to maximally communicate and store information, and to allow the most effective decoding. In machine learning and statistics, by contrast, it is often the case that nature provides a fixed data distribution, and it is the learner’s or statistician’s goal to recover information about this (unknown) distribution. Our goal will be to show how an information theoretic perspective can provide clean answers about and techniques to perform this recovery.

The discovery of fundamental limits forms a central aspect of information theory: the development of results that demonstrate that certain procedures are optimal. Thus, information theoretic tools allow a characterization of the attainable results in a variety of communication and statistical settings. As we explore in the coming chapters in the context of statistical, inferential, and machine learning tasks, this allows us to develop procedures whose optimality we can certify—no better procedure is possible. Such results are useful for a myriad of reasons; we would like to avoid making bad decisions or false inferences, we may realize a task is impossible, and we can explicitly calculate the amount of data necessary for solving different statistical problems.

1.1 Information theory

Information theory focuses on a plethora of deep questions: What is information? How much information content do various signals and data hold? How much information can be reliably transmitted over a noisy communication channel? We will leave delineation of the discipline and answers to these questions to information theorists, instead grossly oversimplifying information theory into two main inquiries, with corresponding chains of tasks.

1. How much information does a signal contain?
2. How much information can a noisy channel reliably transmit?

In this context, we provide two main high-level examples, one for each of these tasks.

Example 1.1.1 (Source coding): The source coding, or data compression problem, is to take information from a source, compress it, decompress it, and recover the original message.

Graphically, we have

$$\text{Source} \rightarrow \text{Compressor} \rightarrow \text{Decompressor} \rightarrow \text{Receiver}$$

The question, then, is how to design a compressor (encoder) and decompressor (decoder) that uses the fewest number of bits to describe a source (or a message) while preserving all the information, in the sense that the receiver receives the correct message with high probability. This fewest number of bits is then the information content of the source (signal). \diamond

Example 1.1.2: The channel coding, or data transmission problem, is the same as the source coding problem of Example 1.1.1, except that between the compressor and decompressor is a source of noise, a *channel*. The graphical representation becomes

$$\text{Source} \rightarrow \text{Compressor} \rightarrow \text{Channel} \rightarrow \text{Decompressor} \rightarrow \text{Receiver}$$

Here we investigate the maximum number of bits that may be sent per each channel use in the sense that the receiver can reconstruct the desired message with low probability of error. Because the channel introduces noise, we require some redundancy, and information theory studies the exact amount of redundancy—in the form of additional bits—that must be sent to allow such reconstruction. \diamond

1.2 Moving to statistics and machine learning

We advocate a study of statistics and machine learning that—broadly—keeps in mind the same views. Let us attempt, then, to shoehorn statistics and machine learning into such source coding and a channel coding problems, which will help to illuminate the perspective that information-theoretic techniques give.

In the analogy with source coding, we observe a sequence of data points X_1, \dots, X_n drawn from some (unknown) distribution P on a space \mathcal{X} . For example, we might be observing species that biologists collect. Then by analogy, we construct a model (often a generative model) that encodes the data using relatively few bits: that is,

$$\text{Source } (P) \xrightarrow{X_1, \dots, X_n} \text{Compressor} \xrightarrow{\hat{P}} \text{Decompressor} \rightarrow \text{Receiver}.$$

Here, we estimate \hat{P} —an empirical version of the distribution P that is easier to describe than the original signal X_1, \dots, X_n —with the hope that we learn information about the generating distribution P , or at least describe it efficiently.

In our analogy with channel coding we can connect to estimation and inference. Consider a statistical problem in which there exists some unknown function f on a space \mathcal{X} that we wish to estimate, and we are able to observe a noisy version of $f(X_i)$ for a series of X_i drawn from a distribution P . Recalling the graphical description of Example 1.1.2, we now have a channel $P(Y | f(X))$ that gives us noisy observations of $f(X)$ for each X_i , but we generally no longer choose the encoder/compressor. That is, we have

$$\text{Source } (P) \xrightarrow{X_1, \dots, X_n} \text{Compressor} \xrightarrow{f(X_1), \dots, f(X_n)} \text{Channel } P(Y | f(X)) \xrightarrow{Y_1, \dots, Y_n} \text{Decompressor}.$$

The estimation—decompression—problem is to either estimate f , or, in some cases, to estimate other aspects of the source probability distribution P . In statistical problems, we do not have

any choice in the design of the compressor f that transforms the original signal X_1, \dots, X_n , which makes it somewhat different from traditional ideas in information theory. In some cases that we explore later—such as experimental design, randomized controlled trials, reinforcement learning and bandits (and associated exploration/exploitation tradeoffs)—we are also able to influence the compression part of the above scheme.

Example 1.2.1: A classical example of the statistical paradigm in this lens is the usual linear regression problem. Here the data X_i belong to \mathbb{R}^d , and the compression function $f(x) = \theta^\top x$ for some vector $\theta \in \mathbb{R}^d$. Then the channel is often of the form

$$Y_i = \underbrace{\theta^\top X_i}_{\text{signal}} + \underbrace{\varepsilon_i}_{\text{noise}},$$

where $\varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$ are independent mean zero normal perturbations. Given a sequence of pairs (X_i, Y_i) , we wish to recover the true θ in the linear model.

In *active learning* or *active sensing* scenarios, also known as (sequential) experimental design, we may choose the sequence X_i so as to better explore properties of θ . As one concrete idea, if we allow infinite *power*, which in this context corresponds to letting $\|X_i\| \rightarrow \infty$ —choosing very “large” vectors x_i —then the signal of $\theta^\top X_i$ should swamp any noise and make estimation easier. \diamond

The remainder of book explores these ideas.

1.3 Outline and chapter discussion

I divide the book into four distinct parts, each of course interacting with the others, but it is possible to read each as a reasonably self-contained unit. The book begins with a review (Chapter 2) that introduces the basic information-theoretic quantities that we discuss: mutual information, entropy, and divergence measures. It is required reading for all the chapters that follow. Chapter 3 provides an overview of exponential family models, which form a core tool in the statistical learning toolbox. Readers familiar with this material, perhaps via a course on generalized linear models, can certainly skip this, but it provides a useful grounding for examples and applications in the subsequent chapters, and so we will dip back into it throughout the book.

Part I of the book covers what I term “stability” based results. At a high level, this means that we ask what can be gained by considering situations where individual observations in a sequence of random variables X_1, \dots, X_n have little effect on various functions of the sequence. We begin in Chapter 4 with concentration inequalities, discussing how sums and related quantities can converge quickly; while this material is essential for the remainder of the chapters, it does not depend on particular information-theoretic techniques. We discuss some heuristic applications to problems in statistical learning—empirical risk minimization—in this section of the book, with Chapter 5 providing results on uniform concentration, with applications to both “generalization”—the standard theoretical tool in machine learning, most typically applying to the accuracy of prediction models—and to estimation problems, which provide various guarantees on estimation of model parameters, which constitute core statistical problems and techniques.

We then turn in Chapter 6 to carefully investigate generalization and convergence guarantees—arguing that functions of a sample X_1, \dots, X_n are representative of the full population P from which the sample is drawn—based on controlling different information-theoretic quantities. In this

context, we develop PAC-Bayesian bounds, and we also use the same framework to present tools to control generalization and convergence in *interactive* data analyses. These types of analyses reflect modern statistics, where one performs some type of data exploration before committing to a fuller analysis, but which breaks classical statistical approaches, because the analysis now depends on the sample. We provide a treatment of more advanced ideas in Chapter 7, where we develop more sophisticated concentration results, such as on random matrices, using core ideas from information theory, which allow us to connect divergence measures to different random processes. Finally, we provide a chapter (Chapter 8) on disclosure limitation and privacy techniques, all of which repose on different notions of stability in distribution.

Part II studies fundamental limits, using information-theoretic techniques to derive *lower bounds* on the possible rates of convergence for various estimation, learning, and other statistical problems. Chapter 9 kicks things off by developing the three major methods for lower bounds: the Assouad, Fano, and Le Cam methods. This chapter shows the basic techniques from which all the other lower bound ideas follow. At a high level, we might consider it, along with Part I, as exhibiting the entire object of study of this book: how do distributions get close to one another, and how can we leverage that closeness? We give a brief treatment of some lower bounding techniques beyond these approaches in Chapter 10, including applications to certain nonparametric problems, as well as a few results that move beyond the typical lower bounds, which apply in expectation, to some that mimic “strong converses” in information theory, meaning that with exceedingly high probability, one cannot hope to achieve anything better than average case error guarantees.

In modern statistical learning problems, one frequently has concerns beyond just statistical risk, such as communication or computational cost, or the privacy of study participants. Accordingly, we develop some of the recent techniques for such problems in Chapter 11 on problems where we wish to obtain optimality guarantees simultaneously along many dimensions, connecting to communication complexity ideas from information theory. Chapter 12 provides a bit of a throwback to estimation with squared error—the most common error metric—introducing the classical statistical tools we have, but shows a few of the more modern applications of the ideas, which re-appear with some frequency. Finally, we conclude the discussion of fundamental limits by looking at testing problems and functional estimation, where one wishes to only estimate a single parameter of a larger model (Chapter 13). While estimating a single scalar might seem, *a priori*, to be simpler than other problems, adequately addressing its complexity requires a fairly nuanced treatment and the introduction of careful information-theoretic tools.

Part III revisits all of our information theoretic notions from Chapter 2, but instead of simply giving definitions and a few consequences, provides operational interpretations of the different information-theoretic quantities, such as entropy. Of course this includes Shannon’s original results on the relationship between coding and entropy (which we cover in the overview Chapter 2.4.1 on information theory), but we also provide an interpretation of entropy and information as measures of uncertainty in statistical experiments and statistical learning, which is a perspective typically missing from information-theoretic treatments of entropy (Chapter 14). Our treatment shows a deep connection between entropy and loss functions used for prediction, where a particular duality allows moving back and forth between them.

We connect these ideas to the problem of *calibration* in Chapter 15, where we ask that a prediction model be valid in that, e.g., on 75% of the days the model provides a prediction of 75% of rain, it rains. We are also able to use these information-theoretic notions of risk, entropy, and losses to connect to problems in optimization and machine learning. In particular, Chapter 16 explores the ways that, if instead of fitting a model to some “true” loss we use an easier-to-optimize surrogate, we essentially lose nothing. This allows us to delineate when (at least in asymptotic senses) it

is possible to computationally efficiently learn good predictors and design good experiments in statistical machine learning problems. Because of the connections with optimization and convex duality, these chapters repose on a nontrivial foundation of convex analysis; we include Appendices (Appendix B and C) that provide a fairly comprehensive review of the results we require. For readers unfamiliar with convex optimization and analysis, I will be the first to admit that these chapters may be tough going—accordingly, we attempt to delineate the big-picture ideas from the nitty-gritty technical conditions necessary for the most general results.

Part IV finishes the book with a treatment of stochastic optimization, online game playing, and minimax problems. Our approach in Chapter 17 takes a modern perspective on stochastic optimization as minimizing random models of functions, and it includes the “book” proofs of convergence of the workhorses of modern machine learning optimization. It also leverages the earlier results on fundamental limits to develop optimality theory for convex optimization in the same framework. Chapter 18 explores online decision-making problems and, more broadly, problems that require exploration and exploitation. This includes bandit problems and some basic questions in causal estimation, where information-theoretic tools allow a clean treatment. The concluding Chapter 19 revisits Chapter 14 on loss functions and predictions, but considers it more in the context of particular games between nature and a statistician/learner. Once again leveraging the perspective on entropy and loss functions we have developed, we are able to provide a generalization of the celebrated redundancy/capacity theorem from information theory, but recast as a game of loss minimization against a nature.

1.4 A remark about measure theory

As this book focuses on a number of fundamental questions in statistics, machine learning, and information theory, fully general statements of the results often require measure theory. Thus, formulae such as $\int f(x)dP(x)$ or $\int f(x)d\mu(x)$ appear. While knowledge of measure theory is certainly useful and may help appreciate the results, it is completely inessential to developing the intuition and, I hope, understanding the proofs and main results. Indeed, the best strategy (for a reader unfamiliar with measure theory) is to simply replace every instance of a formula such as $d\mu(x)$ with dx . The most frequent cases we encounter will be the following: we wish to compute the expectation of a function f of random variable X following distribution P , that is, $\mathbb{E}_P[f(X)]$. Normally, we would write $\mathbb{E}_P[f(X)] = \int f(x)dP(x)$, or sometimes $\mathbb{E}_P[f(X)] = \int f(x)p(x)d\mu(x)$, saying that “ P has density p with respect to the underlying measure μ .” Instead, one may simply (and intuitively) assume that x really has density p over the reals, and instead of computing the integral

$$\mathbb{E}_P[f(X)] = \int f(x)dP(x) \quad \text{or} \quad \mathbb{E}_P[f(X)] = \int f(x)p(x)d\mu(x),$$

assume we may write

$$\mathbb{E}_P[f(X)] = \int f(x)p(x)dx.$$

Nothing will be lost.

Chapter 2

An information theory review

In this first introductory chapter, we discuss and review many of the basic concepts of information theory in effort to introduce them to readers unfamiliar with the tools. Our presentation is relatively brisk, as our main goal is to get to the meat of the chapters on applications of the inequalities and tools we develop, but these provide the starting point for everything in the sequel. One of the main uses of information theory is to prove what, in an information theorist’s lexicon, are known as *converse results*: fundamental limits that guarantee no procedure can improve over a particular benchmark or baseline. We will give the first of these here to preview more of what is to come, as these fundamental limits form one of the core connections between statistics and information theory. The tools of information theory, in addition to their mathematical elegance, also come with strong operational interpretations: they give quite precise answers and explanations for a variety of real engineering and statistical phenomena. We will touch on one of these here (the connection between source coding, or lossless compression, and the Shannon entropy), and much of the remainder of the book will explore more.

2.1 Basics of Information Theory

In this section, we review the basic definitions in information theory, including (Shannon) entropy, KL-divergence, mutual information, and their conditional versions. Before beginning, I must make an apology to any information theorist reading these notes: any time we use a log, it will always be base- e . This is more convenient for our analyses, and it also (later) makes taking derivatives much nicer.

In this first section, we will assume that all distributions are discrete; this makes the quantities somewhat easier to manipulate and allows us to completely avoid any complicated measure-theoretic quantities. In Section 2.2 of this note, we show how to extend the important definitions (for our purposes)—those of KL-divergence and mutual information—to general distributions, where basic ideas such as entropy no longer make sense. However, even in this general setting, we will see we essentially lose no generality by assuming all variables are discrete.

2.1.1 Definitions

Here, we provide the basic definitions of entropy, information, and divergence, assuming the random variables of interest are discrete or have densities with respect to Lebesgue measure.

Entropy: We begin with a central concept in information theory: the entropy. Let P be a distribution on a finite (or countable) set \mathcal{X} , and let p denote the probability mass function associated with P . That is, if X is a random variable distributed according to P , then $P(X = x) = p(x)$. The *entropy of X* (or of P) is defined as

$$H(X) := - \sum_x p(x) \log p(x).$$

Because $p(x) \leq 1$ for all x , it is clear that this quantity is positive. We will show later that if \mathcal{X} is finite, the maximum entropy distribution on \mathcal{X} is the uniform distribution, setting $p(x) = 1/|\mathcal{X}|$ for all x , which has entropy $\log(|\mathcal{X}|)$.

Later in the class, we provide a number of operational interpretations of the entropy. The most common interpretation—which forms the beginning of Shannon’s classical information theory [167]—is via the source-coding theorem. We present Shannon’s source coding theorem in Section 2.4.1, where we show that if we wish to encode a random variable X , distributed according to P , with a k -ary string (i.e. each entry of the string takes on one of k values), then the minimal expected length of the encoding is given by $H(X) = - \sum_x p(x) \log_k p(x)$. Moreover, this is achievable (to within a length of at most 1 symbol) by using Huffman codes (among many other types of codes). As an example of this interpretation, we may consider encoding a random variable X with equi-probable distribution on m items, which has $H(X) = \log(m)$. In base-2, this makes sense: we simply assign an integer to each item and encode each integer with the natural (binary) integer encoding of length $\lceil \log m \rceil$.

We can also define the *conditional entropy*, which is the amount of information left in a random variable after observing another. In particular, we define

$$H(X | Y = y) = - \sum_x p(x | y) \log p(x | y) \quad \text{and} \quad H(X | Y) = \sum_y p(y) H(X | Y = y),$$

where $p(x | y)$ is the p.m.f. of X given that $Y = y$.

Let us now provide a few examples of the entropy of various discrete random variables

Example 2.1.1 (Uniform random variables): As we noted earlier, if a random variable X is uniform on a set of size m , then $H(X) = \log m$. \diamond

Example 2.1.2 (Bernoulli random variables): Let $h_2(p) = -p \log p - (1-p) \log(1-p)$ denote the binary entropy, which is the entropy of a $\text{Bernoulli}(p)$ random variable. \diamond

Example 2.1.3 (Geometric random variables): A random variable X is $\text{Geometric}(p)$, for some $p \in [0, 1]$, if it is supported on $\{1, 2, \dots\}$, and $P(X = k) = (1-p)^{k-1}p$; this is the probability distribution of the number X of $\text{Bernoulli}(p)$ trials until a single success. The entropy of such a random variable is

$$H(X) = - \sum_{k=1}^{\infty} (1-p)^{k-1} p [(k-1) \log(1-p) + \log p] = - \sum_{k=0}^{\infty} (1-p)^k p [k \log(1-p) + \log p].$$

As $\sum_{k=0}^{\infty} \alpha^k = \frac{1}{1-\alpha}$ and $\frac{d}{d\alpha} \frac{1}{1-\alpha} = \frac{1}{(1-\alpha)^2} = \sum_{k=1}^{\infty} k \alpha^{k-1}$, we have

$$H(X) = -p \log(1-p) \cdot \sum_{k=1}^{\infty} k (1-p)^k - p \log p \cdot \sum_{k=1}^{\infty} (1-p)^k = -\frac{1-p}{p} \log(1-p) - (1-p) \log p.$$

As $p \downarrow 0$, we see that $H(X) \uparrow \infty$. \diamond

Example 2.1.4 (A random variable with infinite entropy): While most “reasonable” discrete random variables have finite entropy, it is possible to construct distributions with infinite entropy. Indeed, let X have p.m.f. on $\{2, 3, \dots\}$ defined by

$$p(k) = \frac{A}{k \log^2 k} \quad \text{where} \quad A^{-1} = \sum_{k=2}^{\infty} \frac{1}{k \log^2 k} < \infty,$$

the last sum finite as $\int_2^{\infty} \frac{1}{x \log^{\alpha} x} dx < \infty$ if and only if $\alpha > 1$: for $\alpha = 1$, we have $\int_e^x \frac{1}{t \log t} = \log \log x$, while for $\alpha > 1$, we have

$$\frac{d}{dx} (\log x)^{1-\alpha} = (1-\alpha) \frac{1}{x \log^{\alpha} x}$$

so that $\int_e^{\infty} \frac{1}{t \log^{\alpha} t} dt = \frac{1}{e(1-\alpha)}$. To see that the entropy is infinite, note that

$$H(X) = A \sum_{k \geq 2} \frac{\log A + \log k + 2 \log \log k}{k \log^2 k} \geq A \sum_{k \geq 2} \frac{\log k}{k \log^2 k} - C = \infty,$$

where C is a numerical constant. \diamond

KL-divergence: Now we define two additional quantities, which are actually *much more* fundamental than entropy: they can always be defined for any distributions and any random variables, as they measure distance between distributions. Entropy simply makes no sense for non-discrete random variables, let alone random variables with continuous and discrete components, though it proves useful for some of our arguments and interpretations.

Before defining these quantities, we recall the definition of a convex function $f : \mathbb{R}^k \rightarrow \mathbb{R}$ as any bowl-shaped function, that is, one satisfying

$$f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y) \quad (2.1.1)$$

for all $\lambda \in [0, 1]$, all x, y . The function f is *strictly* convex if the convexity inequality (2.1.1) is strict for $\lambda \in (0, 1)$ and $x \neq y$. We recall a standard result:

Proposition 2.1.5 (Jensen’s inequality). *Let f be convex. Then for any random variable X ,*

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)].$$

Moreover, if f is strictly convex, then $f(\mathbb{E}[X]) < \mathbb{E}[f(X)]$ unless X is constant.

Now we may define and provide a few properties of the KL-divergence. Let P and Q be distributions defined on a discrete set \mathcal{X} . The *KL-divergence* between them is

$$D_{\text{kl}}(P\|Q) := \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)}.$$

We observe immediately that $D_{\text{kl}}(P\|Q) \geq 0$. To see this, we apply Jensen’s inequality (Proposition 2.1.5) to the function $-\log$ and the random variable $q(X)/p(X)$, where X is distributed according to P :

$$\begin{aligned} D_{\text{kl}}(P\|Q) &= -\mathbb{E} \left[\log \frac{q(X)}{p(X)} \right] \geq -\log \mathbb{E} \left[\frac{q(X)}{p(X)} \right] \\ &= -\log \left(\sum_x p(x) \frac{q(x)}{p(x)} \right) = -\log(1) = 0. \end{aligned}$$

Moreover, as log is strictly convex, we have $D_{\text{kl}}(P\|Q) > 0$ unless $P = Q$. Another consequence of the positivity of the KL-divergence is that whenever the set \mathcal{X} is finite with cardinality $|\mathcal{X}| < \infty$, for any random variable X supported on \mathcal{X} we have $H(X) \leq \log |\mathcal{X}|$. Indeed, letting $m = |\mathcal{X}|$, Q be the uniform distribution on \mathcal{X} so that $q(x) = \frac{1}{m}$, and X have distribution P on \mathcal{X} , we have

$$0 \leq D_{\text{kl}}(P\|Q) = \sum_x p(x) \log \frac{p(x)}{q(x)} = -H(X) - \sum_x p(x) \log q(x) = -H(X) + \log m, \quad (2.1.2)$$

so that $H(X) \leq \log m$. Thus, the uniform distribution has the highest entropy over all distributions on the set \mathcal{X} .

Mutual information: Having defined KL-divergence, we may now describe the information content between two random variables X and Y . The *mutual information* $I(X; Y)$ between X and Y is the KL-divergence between their joint distribution and their products (marginal) distributions. More mathematically,

$$I(X; Y) := \sum_{x,y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}. \quad (2.1.3)$$

We can rewrite this in several ways. First, using Bayes' rule, we have $p(x, y)/p(y) = p(x | y)$, so

$$\begin{aligned} I(X; Y) &= \sum_{x,y} p(y)p(x | y) \log \frac{p(x | y)}{p(x)} \\ &= - \sum_x \sum_y p(y)p(x | y) \log p(x) + \sum_y p(y) \sum_x p(x | y) \log p(x | y) \\ &= H(X) - H(X | Y). \end{aligned}$$

Similarly, we have $I(X; Y) = H(Y) - H(Y | X)$, so mutual information can be thought of as the amount of entropy removed (on average) in X by observing Y . We may also think of mutual information as measuring the similarity between the joint distribution of X and Y and their distribution when they are treated as independent.

Comparing the definition (2.1.3) to that for KL-divergence, we see that if P_{XY} is the joint distribution of X and Y , while P_X and P_Y are their marginal distributions (distributions when X and Y are treated independently), then

$$I(X; Y) = D_{\text{kl}}(P_{XY}\|P_X \times P_Y) \geq 0.$$

Moreover, we have $I(X; Y) > 0$ unless X and Y are independent.

As with entropy, we may also define the *conditional information between X and Y given Z* , which is the mutual information between X and Y when Z is observed (on average). That is,

$$I(X; Y | Z) := \sum_z I(X; Y | Z = z)p(z) = H(X | Z) - H(X | Y, Z) = H(Y | Z) - H(Y | X, Z).$$

Entropies of continuous random variables For continuous random variables, we may define an analogue of the entropy known as *differential entropy*, which for a random variable X with density p is defined by

$$h(X) := - \int p(x) \log p(x) dx. \quad (2.1.4)$$

Note that the differential entropy may be negative—it is no longer directly a measure of the number of bits required to describe a random variable X (on average), as was the case for the entropy. We can similarly define the conditional entropy

$$h(X | Y) = - \int p(y) \int p(x | y) \log p(x | y) dx dy.$$

We remark that the conditional differential entropy of X given Y for Y with arbitrary distribution—so long as X has a density—is

$$h(X | Y) = \mathbb{E} \left[- \int p(x | Y) \log p(x | Y) dx \right],$$

where $p(x | y)$ denotes the conditional density of X when $Y = y$. The KL divergence between distributions P and Q with densities p and q becomes

$$D_{\text{kl}}(P \| Q) = \int p(x) \log \frac{p(x)}{q(x)} dx,$$

and similarly, we have the analogues of mutual information as

$$I(X; Y) = \int p(x, y) \log \frac{p(x, y)}{p(x)p(y)} dx dy = h(X) - h(X | Y) = h(Y) - h(Y | X).$$

As we show in the next subsection, we can define the KL-divergence between arbitrary distributions (and mutual information between arbitrary random variables) more generally without requiring discrete or continuous distributions. Before investigating these issues, however, we present a few examples. We also see immediately that for X uniform on a set $[a, b]$, we have $h(X) = \log(b - a)$.

Example 2.1.6 (Entropy of normal random variables): The differential entropy (2.1.4) of a normal random variable is straightforward to compute. Indeed, for $X \sim \mathcal{N}(\mu, \sigma^2)$ we have $p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x - \mu)^2)$, so that

$$h(X) = - \int p(x) \left[\frac{1}{2} \log \frac{1}{2\pi\sigma^2} - \frac{1}{2\sigma^2}(x - \mu)^2 \right] = \frac{1}{2} \log(2\pi\sigma^2) + \frac{\mathbb{E}[(X - \mu)^2]}{2\sigma^2} = \frac{1}{2} \log(2\pi e\sigma^2).$$

For a general multivariate Gaussian, where $X \sim \mathcal{N}(\mu, \Sigma)$ for a vector $\mu \in \mathbb{R}^n$ and $\Sigma \succ 0$ with density $p(x) = \frac{1}{(2\pi)^{n/2} \sqrt{\det(\Sigma)}} \exp(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu))$, we similarly have

$$\begin{aligned} h(X) &= \frac{1}{2} \mathbb{E} \left[n \log(2\pi) + \log \det(\Sigma) + (X - \mu)^\top \Sigma^{-1}(X - \mu) \right] \\ &= \frac{n}{2} \log(2\pi) + \frac{1}{2} \log \det(\Sigma) + \frac{1}{2} \text{tr}(\Sigma \Sigma^{-1}) = \frac{n}{2} \log(2\pi e) + \frac{1}{2} \log \det(e\Sigma). \end{aligned}$$

◇

Continuing our examples with normal distributions, we may compute the divergence between two multivariate Gaussian distributions:

Example 2.1.7 (Divergence between Gaussian distributions): Let P be the multivariate normal $\mathcal{N}(\mu_1, \Sigma)$, and Q be the multivariate normal distribution with mean μ_2 and identical covariance $\Sigma \succ 0$. Then we have that

$$D_{\text{kl}}(P \| Q) = \frac{1}{2} (\mu_1 - \mu_2)^\top \Sigma^{-1} (\mu_1 - \mu_2). \quad (2.1.5)$$

We leave the computation of the identity (2.1.5) to the reader. ◇

An interesting consequence of Example 2.1.7 is that if a random vector X has a given covariance $\Sigma \in \mathbb{R}^{n \times n}$, then the multivariate Gaussian with identical covariance has larger differential entropy. Put another way, differential entropy for random variables with second moments is always maximized by the Gaussian distribution.

Proposition 2.1.8. *Let X be a random vector on \mathbb{R}^n with a density, and assume that $\text{Cov}(X) = \Sigma$. Then for $Z \sim \mathcal{N}(0, \Sigma)$, we have*

$$h(X) \leq h(Z).$$

Proof Without loss of generality, we assume that X has mean 0. Let P be the distribution of X with density p , and let Q be multivariate normal with mean 0 and covariance Σ ; let Z be this random variable. Then

$$\begin{aligned} D_{\text{kl}}(P\|Q) &= \int p(x) \log \frac{p(x)}{q(x)} dx = -h(X) + \int p(x) \left[\frac{n}{2} \log(2\pi) - \frac{1}{2} x^\top \Sigma^{-1} x \right] dx \\ &= -h(X) + h(Z), \end{aligned}$$

because Z has the same covariance as X . As $0 \leq D_{\text{kl}}(P\|Q)$, we have $h(Z) \geq h(X)$ as desired. \square

We remark in passing that the fact that Gaussian random variables have the largest entropy has been used to prove stronger variants of the central limit theorem; see the original results of Barron [16], as well as later quantitative results on the increase of entropy of normalized sums by Artstein et al. [9] and Madiman and Barron [143].

2.1.2 Chain rules and related properties

We now illustrate several of the properties of entropy, KL divergence, and mutual information; these allow easier calculations and analysis.

Chain rules: We begin by describing relationships between collections of random variables X_1, \dots, X_n and individual members of the collection. (Throughout, we use the notation $X_i^j = (X_i, X_{i+1}, \dots, X_j)$ to denote the sequence of random variables from indices i through j .)

For the entropy, we have the simplest chain rule:

$$H(X_1, \dots, X_n) = H(X_1) + H(X_2 | X_1) + \dots + H(X_n | X_1^{n-1}).$$

This follows from the standard decomposition of a probability distribution $p(x, y) = p(x)p(y | x)$. to see the chain rule, then, note that

$$\begin{aligned} H(X, Y) &= - \sum_{x, y} p(x)p(y | x) \log p(x)p(y | x) \\ &= - \sum_x p(x) \sum_y p(y | x) \log p(x) - \sum_x p(x) \sum_y p(y | x) \log p(y | x) = H(X) + H(Y | X). \end{aligned}$$

Now set $X = X_1^{n-1}$, $Y = X_n$, and simply induct.

A related corollary of the definitions of mutual information is the well-known result that *conditioning reduces entropy*:

$$H(X | Y) \leq H(X) \quad \text{because} \quad I(X; Y) = H(X) - H(X | Y) \geq 0.$$

So on average, knowing about a variable Y can only decrease your uncertainty about X . That conditioning reduces entropy for continuous random variables is also immediate, as for X continuous we have $I(X; Y) = h(X) - h(X | Y) \geq 0$, so that $h(X) \geq h(X | Y)$.

Chain rules for information and divergence: As another immediate corollary to the chain rule for entropy, we see that mutual information also obeys a chain rule:

$$I(X; Y_1^n) = \sum_{i=1}^n I(X; Y_i \mid Y_1^{i-1}).$$

Indeed, we have

$$I(X; Y_1^n) = H(Y_1^n) - H(Y_1^n \mid X) = \sum_{i=1}^n [H(Y_i \mid Y_1^{i-1}) - H(Y_i \mid X, Y_1^{i-1})] = \sum_{i=1}^n I(X; Y_i \mid Y_1^{i-1}).$$

The KL-divergence obeys similar chain rules, making mutual information and KL-divergence measures useful tools for evaluation of distances and relationships between groups of random variables.

As a second example, suppose that the distribution $P = P_1 \times P_2 \times \cdots \times P_n$, and $Q = Q_1 \times \cdots \times Q_n$, that is, that P and Q are product distributions over independent random variables $X_i \sim P_i$ or $X_i \sim Q_i$. Then we immediately have the tensorization identity

$$D_{\text{kl}}(P \parallel Q) = D_{\text{kl}}(P_1 \times \cdots \times P_n \parallel Q_1 \times \cdots \times Q_n) = \sum_{i=1}^n D_{\text{kl}}(P_i \parallel Q_i).$$

We remark in passing that these two identities hold for arbitrary distributions P_i and Q_i or random variables X, Y . As a final tensorization identity, we consider a more general chain rule for KL-divergences, which will frequently be useful. We abuse notation temporarily, and for random variables X and Y with distributions P and Q , respectively, we denote

$$D_{\text{kl}}(X \parallel Y) := D_{\text{kl}}(P \parallel Q).$$

In analogy to the entropy, we can also define the *conditional KL divergence*. Let X and Y have distributions $P_{X|z}$ and $P_{Y|z}$ conditioned on $Z = z$, respectively. Then we define

$$D_{\text{kl}}(X \parallel Y \mid Z) = \mathbb{E}_Z[D_{\text{kl}}(P_{X|Z} \parallel P_{Y|Z})],$$

so that if Z is discrete we have $D_{\text{kl}}(X \parallel Y \mid Z) = \sum_z p(z) D_{\text{kl}}(P_{X|z} \parallel P_{Y|z})$. With this notation, we have the chain rule

$$D_{\text{kl}}(X_1, \dots, X_n \parallel Y_1, \dots, Y_n) = \sum_{i=1}^n D_{\text{kl}}(X_i \parallel Y_i \mid X_1^{i-1}), \quad (2.1.6)$$

because (in the discrete case, which—as we discuss presently—is fully general for this purpose) for distributions P_{XY} and Q_{XY} we have

$$\begin{aligned} D_{\text{kl}}(P_{XY} \parallel Q_{XY}) &= \sum_{x,y} p(x,y) \log \frac{p(x,y)}{q(x,y)} = \sum_{x,y} p(x)p(y \mid x) \left[\log \frac{p(y \mid x)}{q(y \mid x)} + \log \frac{p(x)}{q(x)} \right] \\ &= \sum_x p(x) \log \frac{p(x)}{q(x)} + \sum_x p(x) \sum_y p(y \mid x) \log \frac{p(y \mid x)}{q(y \mid x)}, \end{aligned}$$

where the final equality uses that $\sum_y p(y \mid x) = 1$ for all x . In different notation, if we let P and Q be any distributions on $\mathcal{X}_1 \times \cdots \times \mathcal{X}_n$, and define $P_i(A \mid x_1^{i-1}) = P(X_i \in A \mid X_1^{i-1} = x_1^{i-1})$, and similarly for Q_i , we have the following:

Lemma 2.1.9. *Let P, Q be distributions on $\mathcal{X}_1 \times \cdots \times \mathcal{X}_n$. Then*

$$D_{\text{kl}}(P \| Q) = \sum_{i=1}^n \mathbb{E}_P[D_{\text{kl}}(P_i(\cdot | X_1^{i-1}) \| Q_i(\cdot | X_1^{i-1}))].$$

Expanding upon this, we give several *tensorization* identities, showing how to transform questions about the joint distribution of many random variables to simpler questions about their marginals. As a first example, we see that as a consequence of the fact that conditioning decreases entropy, we see that for any sequence of (discrete or continuous, as appropriate) random variables, we have

$$H(X_1, \dots, X_n) \leq H(X_1) + \cdots + H(X_n) \quad \text{and} \quad h(X_1, \dots, X_n) \leq h(X_1) + \cdots + h(X_n).$$

Both equalities hold with equality if and only if X_1, \dots, X_n are mutually independent. (The only if follows because $I(X; Y) > 0$ whenever X and Y are not independent, by Jensen's inequality and the fact that $D_{\text{kl}}(P \| Q) > 0$ unless $P = Q$.)

We return to information and divergence now. Suppose that random variables Y_i are independent conditional on X , meaning that

$$P(Y_1 = y_1, \dots, Y_n = y_n | X = x) = P(Y_1 = y_1 | X = x) \cdots P(Y_n = y_n | X = x).$$

Such scenarios are common—as we shall see—when we make multiple observations from a fixed distribution parameterized by some X . Then we have the inequality

$$\begin{aligned} I(X; Y_1, \dots, Y_n) &= \sum_{i=1}^n [H(Y_i | Y_1^{i-1}) - H(Y_i | X, Y_1^{i-1})] \\ &= \sum_{i=1}^n [H(Y_i | Y_1^{i-1}) - H(Y_i | X)] \leq \sum_{i=1}^n [H(Y_i) - H(Y_i | X)] = \sum_{i=1}^n I(X; Y_i), \end{aligned} \tag{2.1.7}$$

where the inequality follows because conditioning reduces entropy.

2.1.3 Data processing inequalities:

A standard problem in information theory (and statistical inference) is to understand the degradation of a signal after it is passed through some noisy channel (or observation process). The simplest of such results, which we will use frequently, is that we can only lose information by adding noise. In particular, assume we have the Markov chain

$$X \rightarrow Y \rightarrow Z.$$

Then we obtain the classical *data processing inequality*.

Proposition 2.1.10. *With the above Markov chain, we have $I(X; Z) \leq I(X; Y)$.*

Proof We expand the mutual information $I(X; Y, Z)$ in two ways:

$$\begin{aligned} I(X; Y, Z) &= I(X; Z) + I(X; Y | Z) \\ &= I(X; Y) + \underbrace{I(X; Z | Y)}_{=0}, \end{aligned}$$

where we note that the final equality follows because X is independent of Z given Y :

$$I(X; Z | Y) = H(X | Y) - H(X | Y, Z) = H(X | Y) - H(X | Y) = 0.$$

Since $I(X; Y | Z) \geq 0$, this gives the result. \square

There are related data processing inequalities for the KL-divergence—which we generalize in the next section—as well. In this case, we may consider a simple Markov chain $X \rightarrow Z$. If we let P_1 and P_2 be distributions on X and Q_1 and Q_2 be the induced distributions on Z , that is, $Q_i(A) = \int \mathbb{P}(Z \in A | x) dP_i(x)$, then we have

$$D_{\text{kl}}(Q_1 \| Q_2) \leq D_{\text{kl}}(P_1 \| P_2),$$

the basic KL-divergence data processing inequality. A consequence of this is that, for any function f and random variables X and Y on the same space, we have

$$D_{\text{kl}}(f(X) \| f(Y)) \leq D_{\text{kl}}(X \| Y).$$

We explore these data processing inequalities more when we generalize KL-divergences in the next section and in the exercises.

2.2 General divergence measures and definitions

Having given our basic definitions of mutual information and divergence, we now show how the definitions of KL-divergence and mutual information extend to arbitrary distributions P and Q and arbitrary sets \mathcal{X} . This requires a bit of setup, including defining set algebras (which, we will see, simply correspond to quantization of the set \mathcal{X}), but allows us to define divergences in full generality.

2.2.1 Partitions, algebras, and quantizers

Let \mathcal{X} be an arbitrary space. A *quantizer* on \mathcal{X} is any function that maps \mathcal{X} to a finite collection of integers. That is, fixing $m < \infty$, a quantizer is any function $\mathbf{q} : \mathcal{X} \rightarrow \{1, \dots, m\}$. In particular, a quantizer \mathbf{q} partitions the space \mathcal{X} into the subsets of $x \in \mathcal{X}$ for which $\mathbf{q}(x) = i$. A related notion—we will see the precise relationship presently—is that of an algebra of sets on \mathcal{X} . We say that a collection of sets \mathcal{A} is an *algebra* on \mathcal{X} if the following are true:

1. The set $\mathcal{X} \in \mathcal{A}$.
2. The collection of sets \mathcal{A} is closed under finite set operations: union, intersection, and complementation. That is, $A, B \in \mathcal{A}$ implies that $A^c \in \mathcal{A}$, $A \cap B \in \mathcal{A}$, and $A \cup B \in \mathcal{A}$.

There is a 1-to-1 correspondence between quantizers—and their associated partitions of the set \mathcal{X} —and finite algebras on a set \mathcal{X} , which we discuss briefly.¹ It should be clear that there is a one-to-one correspondence between finite *partitions* of the set \mathcal{X} and quantizers \mathbf{q} , so we must argue that finite partitions of \mathcal{X} are in one-to-one correspondence with finite algebras defined over \mathcal{X} .

¹Pedantically, this one-to-one correspondence holds up to permutations of the partition induced by the quantizer.

In one direction, we may consider a quantizer $\mathbf{q} : \mathcal{X} \rightarrow \{1, \dots, m\}$. Let the sets A_1, \dots, A_m be the partition associated with \mathbf{q} , that is, for $x \in A_i$ we have $\mathbf{q}(x) = i$, or $A_i = \mathbf{q}^{-1}(\{i\})$. Then we may define an algebra $\mathcal{A}_{\mathbf{q}}$ as the collection of all finite set operations performed on A_1, \dots, A_m (note that this is a finite collection, as finite set operations performed on the partition A_1, \dots, A_m induce only a finite collection of sets).

For the other direction, consider a finite algebra \mathcal{A} over the set \mathcal{X} . We can then construct a quantizer $\mathbf{q}_{\mathcal{A}}$ that corresponds to this algebra. To do so, we define an *atom* of \mathcal{A} as any non-empty set $A \in \mathcal{A}$ such that if $B \subset A$ and $B \in \mathcal{A}$, then $B = A$ or $B = \emptyset$. That is, the atoms of \mathcal{A} are the “smallest” sets in \mathcal{A} . We claim there is a unique partition of \mathcal{X} with atomic sets from \mathcal{A} ; we prove this inductively.

Base case: There is at least 1 atomic set, as \mathcal{A} is finite; call it A_1 .

Induction step: Assume we have atomic sets $A_1, \dots, A_k \in \mathcal{A}$. Let $B = (A_1 \cup \dots \cup A_k)^c$ be their complement, which we assume is non-empty (otherwise we have a partition of \mathcal{X} into atomic sets). The complement B is either atomic, in which case the sets $\{A_1, A_2, \dots, A_k, B\}$ are a partition of \mathcal{X} consisting of atoms of \mathcal{A} , or B is not atomic. If B is not atomic, consider all the sets of the form $A \cap B$ for $A \in \mathcal{A}$. Each of these belongs to \mathcal{A} , and at least one of them is atomic, as there is a finite number of them. This means there is a non-empty set $A_{k+1} \subset B$ such that A_{k+1} is atomic.

By repeating this induction, which must stop at some finite index m as \mathcal{A} is finite, we construct a collection A_1, \dots, A_m of disjoint atomic sets in \mathcal{A} for which $\cup_i A_i = \mathcal{X}$. (The uniqueness is an exercise for the reader.) Thus we may define the quantizer $\mathbf{q}_{\mathcal{A}}$ via

$$\mathbf{q}_{\mathcal{A}}(x) = i \quad \text{when } x \in A_i.$$

2.2.2 KL-divergence

In this section, we present the general definition of a KL-divergence, which holds for *any* pair of distributions. Let P and Q be distributions on a space \mathcal{X} . Now, let \mathcal{A} be a finite algebra on \mathcal{X} (as in the previous section, this is equivalent to picking a partition of \mathcal{X} and then constructing the associated algebra), and assume that its atoms are $\text{atoms}(\mathcal{A})$. The KL-divergence between P and Q *conditioned on* \mathcal{A} is

$$D_{\text{kl}}(P \| Q \mid \mathcal{A}) := \sum_{A \in \text{atoms}(\mathcal{A})} P(A) \log \frac{P(A)}{Q(A)}.$$

That is, we simply sum over the partition of \mathcal{X} . Another way to write this is as follows. Let $\mathbf{q} : \mathcal{X} \rightarrow \{1, \dots, m\}$ be a quantizer, and define the sets $A_i = \mathbf{q}^{-1}(\{i\})$ to be the pre-images of each i (i.e. the different quantization regions, or the partition of \mathcal{X} that \mathbf{q} induces). Then the *quantized* KL-divergence between P and Q is

$$D_{\text{kl}}(P \| Q \mid \mathbf{q}) := \sum_{i=1}^m P(A_i) \log \frac{P(A_i)}{Q(A_i)}.$$

We may now give the fully general definition of KL-divergence: the KL-divergence between P and Q is defined as

$$\begin{aligned} D_{\text{kl}}(P \| Q) &:= \sup \{ D_{\text{kl}}(P \| Q \mid \mathcal{A}) \mid \text{such that } \mathcal{A} \text{ is a finite algebra on } \mathcal{X} \} \\ &= \sup \{ D_{\text{kl}}(P \| Q \mid \mathbf{q}) \mid \text{such that } \mathbf{q} \text{ quantizes } \mathcal{X} \}. \end{aligned} \tag{2.2.1}$$

This also gives a rigorous definition of mutual information. Indeed, if X and Y are random variables with joint distribution P_{XY} and marginal distributions P_X and P_Y , we simply define

$$I(X; Y) = D_{\text{kl}}(P_{XY} \| P_X \times P_Y).$$

When P and Q have densities p and q , the definition (2.2.1) reduces to

$$D_{\text{kl}}(P \| Q) = \int_{\mathbb{R}} p(x) \log \frac{p(x)}{q(x)} dx,$$

while if P and Q both have probability mass functions p and q , then—as we see in Exercise 2.6—the definition (2.2.1) is equivalent to

$$D_{\text{kl}}(P \| Q) = \sum_x p(x) \log \frac{p(x)}{q(x)},$$

precisely as in the discrete case.

We remark in passing that if the set \mathcal{X} is a product space, meaning that $\mathcal{X} = \mathcal{X}_1 \times \mathcal{X}_2 \times \cdots \times \mathcal{X}_n$ for some $n < \infty$ (this is the case for mutual information, for example), then we may assume our quantizer *always* quantizes sets of the form $A = A_1 \times A_2 \times \cdots \times A_n$, that is, Cartesian products. Written differently, when we consider algebras on \mathcal{X} , the atoms of the algebra may be assumed to be Cartesian products of sets, and our partitions of \mathcal{X} can always be taken as Cartesian products. (See Gray [104, Chapter 5].) Written slightly differently, if P and Q are distributions on $\mathcal{X} = \mathcal{X}_1 \times \cdots \times \mathcal{X}_n$ and \mathbf{q}^i is a quantizer for the set \mathcal{X}_i (inducing the partition $A_1^i, \dots, A_{m_i}^i$ of \mathcal{X}_i) we may define

$$D_{\text{kl}}(P \| Q \mid \mathbf{q}^1, \dots, \mathbf{q}^n) = \sum_{j_1, \dots, j_n} P(A_{j_1}^1 \times A_{j_2}^2 \times \cdots \times A_{j_n}^n) \log \frac{P(A_{j_1}^1 \times A_{j_2}^2 \times \cdots \times A_{j_n}^n)}{Q(A_{j_1}^1 \times A_{j_2}^2 \times \cdots \times A_{j_n}^n)}.$$

Then the general definition (2.2.1) of KL-divergence specializes to

$$D_{\text{kl}}(P \| Q) = \sup \{ D_{\text{kl}}(P \| Q \mid \mathbf{q}^1, \dots, \mathbf{q}^n) \text{ such that } \mathbf{q}^i \text{ quantizes } \mathcal{X}_i \}.$$

So we only need consider “rectangular” sets in the definitions of KL-divergence.

Measure-theoretic definition of KL-divergence If you have never seen measure theory before, skim this section; while the notation may be somewhat intimidating, it is fine to always consider only continuous or fully discrete distributions. We will describe an interpretation that will mean for our purposes that one never needs to really think about measure theoretic issues.

The general definition (2.2.1) of KL-divergence is equivalent to the following. Let μ be a measure on \mathcal{X} , and assume that P and Q are absolutely continuous with respect to μ , with densities p and q , respectively. (For example, take $\mu = P + Q$.) Then

$$D_{\text{kl}}(P \| Q) = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} d\mu(x). \quad (2.2.2)$$

The proof of this fact is somewhat involved, requiring the technology of Lebesgue integration. (See Gray [104, Chapter 5].)

For those who have not seen measure theory, the interpretation of the equality (2.2.2) should be as follows. When integrating a function $f(x)$, replace $\int f(x) d\mu(x)$ with one of two pairs of symbols: one may simply think of $d\mu(x)$ as dx , so that we are performing standard integration $\int f(x) dx$, or one should think of the integral operation $\int f(x) d\mu(x)$ as summing the argument of the integral, so $d\mu(x) = 1$ and $\int f(x) d\mu(x) = \sum_x f(x)$. (This corresponds to μ being “counting measure” on \mathcal{X} .)

2.2.3 f -divergences

A more general notion of divergence is the so-called f -divergence, or Ali-Silvey divergence [6, 59] (see also the alternate interpretations in the article by Liese and Vajda [137]). Here, the definition is as follows. Let P and Q be probability distributions on the set \mathcal{X} , and let $f : \mathbb{R}_+ \rightarrow \mathbb{R}$ be a convex function satisfying $f(1) = 0$. If \mathcal{X} is a discrete set, then the f -divergence between P and Q is

$$D_f(P\|Q) := \sum_x q(x) f\left(\frac{p(x)}{q(x)}\right).$$

More generally, for any set \mathcal{X} and a quantizer $\mathbf{q} : \mathcal{X} \rightarrow \{1, \dots, m\}$, letting $A_i = \mathbf{q}^{-1}(\{i\}) = \{x \in \mathcal{X} \mid \mathbf{q}(x) = i\}$ be the partition the quantizer induces, we can define the quantized divergence

$$D_f(P\|Q \mid \mathbf{q}) = \sum_{i=1}^m Q(A_i) f\left(\frac{P(A_i)}{Q(A_i)}\right),$$

and the general definition of an f divergence is (in analogy with the definition (2.2.1) of general KL divergences)

$$D_f(P\|Q) := \sup \{D_f(P\|Q \mid \mathbf{q}) \mid \mathbf{q} \text{ quantizes } \mathcal{X}\}. \quad (2.2.3)$$

The definition (2.2.3) shows that, any time we have computations involving f -divergences—such as KL-divergence or mutual information—it is no loss of generality, when performing the computations, to assume that all distributions have finite discrete support. There is a measure-theoretic version of the definition (2.2.3) which is frequently easier to use. Assume w.l.o.g. that P and Q are absolutely continuous with respect to the base measure μ . The f divergence between P and Q is then

$$D_f(P\|Q) := \int_{\mathcal{X}} q(x) f\left(\frac{p(x)}{q(x)}\right) d\mu(x). \quad (2.2.4)$$

This definition, it turns out, is not *quite* as general as we would like—in particular, it is unclear how we should define the integral for points x such that $q(x) = 0$. With that in mind, we recall that the perspective transform (see Appendices B.1.1 and B.3.3) of a function $f : \mathbb{R} \rightarrow \mathbb{R}$ is defined by $\text{pers}(f)(t, u) = uf(t/u)$ if $u > 0$ and by $+\infty$ if $u \leq 0$. This function is convex in its arguments (Proposition B.3.12). In fact, this is not quite enough for the fully correct definition. The *closure* of a convex function f is $\text{cl } f(x) = \sup\{\ell(x) \mid \ell \leq f, \ell \text{ linear}\}$, the supremum over all linear functions that globally lower bound f . Then [111, Proposition IV.2.2.2] the closer of $\text{pers}(f)$ is defined, for any $t' \in \text{int dom } f$, by

$$\text{cl pers}(f)(t, u) = \begin{cases} uf(t/u) & \text{if } u > 0 \\ \lim_{\alpha \downarrow 0} \alpha f(t' - t + t/\alpha) & \text{if } u = 0 \\ +\infty & \text{if } u < 0. \end{cases}$$

(The choice of t' does not affect the definition.) Then the fully general formula expressing the f -divergence is

$$D_f(P\|Q) = \int_{\mathcal{X}} \text{cl pers}(f)(p(x), q(x)) d\mu(x). \quad (2.2.5)$$

This is what we mean by equation (2.2.4), which we use without comment.

In the exercises, we explore several properties of f -divergences, including the quantized representation (2.2.3), showing different data processing inequalities and orderings of quantizers based

on the fineness of their induced partitions. Broadly, f -divergences satisfy essentially the same properties as KL-divergence, such as data-processing inequalities, and they provide a generalization of mutual information. We explore f -divergences from additional perspectives later—they are important both for optimality in estimation and related to consistency and prediction problems, as we discuss in Chapter 16.5.

Examples We give several examples of f -divergences here; in Section 9.2.2 we provide a few examples of their uses as well as providing a few natural inequalities between them.

Example 2.2.1 (KL-divergence): By taking $f(t) = t \log t$, which is convex and satisfies $f(1) = 0$, we obtain $D_f(P\|Q) = D_{\text{kl}}(P\|Q)$. \diamond

Example 2.2.2 (KL-divergence, reversed): By taking $f(t) = -\log t$, we obtain $D_f(P\|Q) = D_{\text{kl}}(Q\|P)$. \diamond

Example 2.2.3 (Total variation distance): The total variation distance between probability distributions P and Q defined on a set \mathcal{X} is the maximum difference between probabilities they assign on subsets of \mathcal{X} :

$$\|P - Q\|_{\text{TV}} := \sup_{A \subset \mathcal{X}} |P(A) - Q(A)| = \sup_{A \subset \mathcal{X}} (P(A) - Q(A)), \quad (2.2.6)$$

where the second equality follows by considering compliments $P(A^c) = 1 - P(A)$. The total variation distance, as we shall see later, is important for verifying the optimality of different tests, and appears in the measurement of difficulty of solving hypothesis testing problems. The choice $f(t) = \frac{1}{2}|t - 1|$, we obtain the total variation distance, that is, $\|P - Q\|_{\text{TV}} = D_f(P\|Q)$. There are several alternative characterizations, which we provide as Lemma 2.2.4 next; it will be useful in the sequel when we develop inequalities relating the divergences. \diamond

Lemma 2.2.4. *Let P, Q be probability measures with densities p, q with respect to a base measure μ and $f(t) = \frac{1}{2}|t - 1|$. Then*

$$\begin{aligned} \|P - Q\|_{\text{TV}} &= D_f(P\|Q) = \frac{1}{2} \int |p(x) - q(x)| d\mu(x) \\ &= \int [p(x) - q(x)]_+ d\mu(x) = \int [q(x) - p(x)]_+ d\mu(x) \\ &= P(dP/dQ > 1) - Q(dP/dQ > 1) = Q(dQ/dP > 1) - P(dQ/dP > 1). \end{aligned}$$

In particular, the set $A = \{x \mid p(x)/q(x) \geq 1\}$ maximizes $P(B) - Q(B)$ over $B \subset \mathcal{X}$ and so achieves $\|P - Q\|_{\text{TV}} = P(A) - Q(A)$.

Proof Eliding the measure-theoretic details,² we immediately have

$$\begin{aligned} D_f(P\|Q) &= \frac{1}{2} \int \left| \frac{p(x)}{q(x)} - 1 \right| q(x) d\mu(x) = \frac{1}{2} \int |p(x) - q(x)| d\mu(x) \\ &= \frac{1}{2} \int_{x:p(x)>q(x)} [p(x) - q(x)] d\mu(x) + \frac{1}{2} \int_{x:q(x)>p(x)} [q(x) - p(x)] d\mu(x) \\ &= \frac{1}{2} \int [p(x) - q(x)]_+ d\mu(x) + \frac{1}{2} \int [q(x) - p(x)]_+ d\mu(x). \end{aligned}$$

²To make this fully rigorous, we would use the Hahn decomposition of the signed measure $P - Q$ to recognize that $\int f(dP - dQ) = \int f[dP - dQ]_+ - \int f[dQ - dP]_+$ for any integrable f .

Considering the last integral $\int [q(x) - p(x)]_+ d\mu(x)$, we see that the set $A = \{x : q(x) > p(x)\}$ satisfies

$$Q(A) - P(A) = \int_A (q(x) - p(x)) d\mu(x) \geq \int_B (q(x) - p(x)) d\mu(x) = Q(B) - P(B)$$

for any set B , as any $x \in B \setminus A$ clearly satisfies $q(x) - p(x) \leq 0$. \square

Example 2.2.5 (Hellinger distance): The *Hellinger distance* between probability distributions P and Q defined on a set \mathcal{X} is generated by the function $f(t) = (\sqrt{t} - 1)^2 = t - 2\sqrt{t} + 1$. The Hellinger distance is then

$$d_{\text{hel}}(P, Q)^2 := \frac{1}{2} \int (\sqrt{p(x)} - \sqrt{q(x)})^2 d\mu(x). \quad (2.2.7)$$

The non-squared version $d_{\text{hel}}(P, Q)$ is indeed a distance between probability measures P and Q . It is sometimes convenient to rewrite the Hellinger distance in terms of the *affinity* between P and Q , as

$$d_{\text{hel}}(P, Q)^2 = \frac{1}{2} \int (p(x) + q(x) - 2\sqrt{p(x)q(x)}) d\mu(x) = 1 - \int \sqrt{p(x)q(x)} d\mu(x), \quad (2.2.8)$$

which makes clear that $d_{\text{hel}}(P, Q) \in [0, 1]$ is on roughly the same scale as the variation distance; we will say more later. \diamond

Example 2.2.6 (χ^2 divergence): The χ^2 -divergence is generated by taking $f(t) = (t - 1)^2$, so that

$$D_{\chi^2}(P \| Q) := \int \left(\frac{p(x)}{q(x)} - 1 \right)^2 q(x) d\mu(x) = \int \frac{p(x)^2}{q(x)} d\mu(x) - 1, \quad (2.2.9)$$

where the equality is immediate because $\int p d\mu = \int q d\mu = 1$. \diamond

2.2.4 Inequalities and relationships between divergences

Important to our development will come will be different families of inequalities relating the different divergence measures. These inequalities will be particularly important because, in some cases, different distributions admit easy calculations with some divergences, such as KL or χ^2 divergence, but it can be challenging to work with others that may be more “natural” for a particular problem. Most importantly, replacing a variation distance by bounding it with an alternative divergence is often convenient for analyzing the properties of product distributions (as will become apparent in Chapter 9). We record several of these results here, making a passing connection to mutual information as well.

The first inequality shows that the Hellinger distance and variation distance roughly generate the same topology on collections of distributions, as they upper and lower bound the other (if we tolerate polynomial losses).

Proposition 2.2.7. *The total variation distance and Hellinger distance satisfy*

$$d_{\text{hel}}^2(P, Q) \leq \|P - Q\|_{\text{TV}} \leq d_{\text{hel}}(P, Q) \sqrt{2 - d_{\text{hel}}^2(P, Q)}.$$

Proof We begin with the upper bound. We have by Hölder's inequality that

$$\begin{aligned} \frac{1}{2} \int |p(x) - q(x)| d\mu(x) &= \int |\sqrt{p(x)} - \sqrt{q(x)}| \cdot |\sqrt{p(x)} + \sqrt{q(x)}| d\mu(x) \\ &\leq \left(\frac{1}{2} \int (\sqrt{p(x)} - \sqrt{q(x)})^2 d\mu(x) \right)^{\frac{1}{2}} \left(\frac{1}{2} \int (\sqrt{p(x)} + \sqrt{q(x)})^2 d\mu(x) \right)^{\frac{1}{2}} \\ &= d_{\text{hel}}(P, Q) \left(1 + \int \sqrt{p(x)q(x)} d\mu(x) \right)^{\frac{1}{2}}. \end{aligned}$$

As in Example 2.2.5, we have $\int \sqrt{p(x)q(x)} d\mu(x) = 1 - d_{\text{hel}}(P, Q)^2$, so this (along with the representation Lemma 2.2.4 for variation distance) implies

$$\|P - Q\|_{\text{TV}} = \frac{1}{2} \int |p(x) - q(x)| d\mu(x) \leq d_{\text{hel}}(P, Q) (2 - d_{\text{hel}}^2(P, Q))^{\frac{1}{2}}.$$

For the lower bound on total variation, note that for any $a, b \in \mathbb{R}_+$, we have $a + b - 2\sqrt{ab} \leq |a - b|$ (check the cases $a > b$ and $a < b$ separately); thus

$$d_{\text{hel}}^2(P, Q) = \frac{1}{2} \int [p(x) + q(x) - 2\sqrt{p(x)q(x)}] d\mu(x) \leq \frac{1}{2} \int |p(x) - q(x)| d\mu(x),$$

as desired. \square

Several important inequalities relate the variation distance to the KL-divergence. We state two important inequalities in the next proposition, both of which are important enough to justify their own names.

Proposition 2.2.8. *The total variation distance satisfies the following relationships.*

(a) Pinsker's inequality: for any distributions P and Q ,

$$\|P - Q\|_{\text{TV}}^2 \leq \frac{1}{2} D_{\text{kl}}(P \| Q). \quad (2.2.10)$$

(b) The Bretagnolle-Huber inequality: for any distributions P and Q ,

$$\|P - Q\|_{\text{TV}} \leq \sqrt{1 - \exp(-D_{\text{kl}}(P \| Q))} \leq 1 - \frac{1}{2} \exp(-D_{\text{kl}}(P \| Q)).$$

Proof Exercise 2.19 outlines one proof of Pinsker's inequality using the data processing inequality (Proposition 2.2.13). We present an alternative via the Cauchy-Schwarz inequality. Using the definition (2.2.1) of the KL-divergence, we may assume without loss of generality that P and Q are finitely supported, say with p.m.f.s p_1, \dots, p_m and q_1, \dots, q_m . Define the negative entropy function $h(p) = \sum_{i=1}^m p_i \log p_i$. Then showing that $D_{\text{kl}}(P \| Q) \geq 2 \|P - Q\|_{\text{TV}}^2 = \frac{1}{2} \|p - q\|_1^2$ is equivalent to showing that

$$h(p) \geq h(q) + \langle \nabla h(q), p - q \rangle + \frac{1}{2} \|p - q\|_1^2, \quad (2.2.11)$$

because by inspection $h(p) - h(q) - \langle \nabla h(q), p - q \rangle = \sum_i p_i \log \frac{p_i}{q_i}$. We do this via a Taylor expansion: we have

$$\nabla h(p) = [\log p_i + 1]_{i=1}^m \quad \text{and} \quad \nabla^2 h(p) = \text{diag}([1/p_i]_{i=1}^m).$$

By Taylor's theorem, there is some $\tilde{p} = (1 - t)p + tq$, where $t \in [0, 1]$, such that

$$h(p) = h(q) + \langle \nabla h(q), p - q \rangle + \frac{1}{2} \langle p - q, \nabla^2 h(\tilde{p})(p - q) \rangle.$$

But looking at the final quadratic, we have for any vector v and any $p \geq 0$ satisfying $\sum_i p_i = 1$,

$$\langle v, \nabla^2 h(\tilde{p})v \rangle = \sum_{i=1}^m \frac{v_i^2}{p_i} = \|p\|_1 \sum_{i=1}^m \frac{v_i^2}{p_i} \geq \left(\sum_{i=1}^m \sqrt{p_i} \frac{|v_i|}{\sqrt{p_i}} \right)^2 = \|v\|_1^2,$$

where the inequality follows from Cauchy-Schwarz applied to the vectors $[\sqrt{p_i}]_i$ and $[|v_i|/\sqrt{p_i}]_i$. Thus inequality (2.2.11) holds.

For the claim (b), we use Proposition 2.2.7. Let $a = \int \sqrt{p(x)q(x)} d\mu(x)$ be a shorthand for the affinity, so that $d_{\text{hel}}^2(P, Q) = 1 - a$. Then Proposition 2.2.7 gives $\|P - Q\|_{\text{TV}} \leq \sqrt{1 - a}\sqrt{1 + a} = \sqrt{1 - a^2}$. Now apply Jensen's inequality to the exponential: we have

$$\begin{aligned} \int \sqrt{p(x)q(x)} d\mu(x) &= \int \sqrt{\frac{q(x)}{p(x)}} p(x) d\mu(x) = \int \exp\left(\frac{1}{2} \log \frac{q(x)}{p(x)}\right) p(x) d\mu(x) \\ &\geq \exp\left(\frac{1}{2} \int p(x) \log \frac{q(x)}{p(x)} d\mu(x)\right) = \exp\left(-\frac{1}{2} D_{\text{kl}}(P\|Q)\right). \end{aligned}$$

In particular, $\sqrt{1 - a^2} \leq \sqrt{1 - \exp(-\frac{1}{2} D_{\text{kl}}(P\|Q))}^2$, which is the first claim of part (b). For the second, note that $\sqrt{1 - c} \leq 1 - \frac{1}{2}c$ for $c \in [0, 1]$ by concavity of the square root. \square

We also have the following bounds on the Hellinger distance in terms of the KL-divergence, and that in terms of the χ^2 -divergence.

Proposition 2.2.9. *For any distributions P, Q ,*

$$2d_{\text{hel}}^2(P, Q) \leq D_{\text{kl}}(P\|Q) \leq \log(1 + D_{\chi^2}(P\|Q)) \leq D_{\chi^2}(P\|Q).$$

Proof For the first inequality, note that $\log x \leq x - 1$ by concavity, or $1 - x \leq -\log x$, so that

$$\begin{aligned} 2d_{\text{hel}}^2(P, Q) &= 2 - 2 \int \sqrt{p(x)q(x)} d\mu(x) \\ &= 2 \int p(x) \left(1 - \sqrt{\frac{q(x)}{p(x)}}\right) d\mu(x) \leq 2 \int p(x) \log \sqrt{\frac{p(x)}{q(x)}} d\mu(x) = D_{\text{kl}}(P\|Q). \end{aligned}$$

The last two inequalities are simple: by Jensen's inequality, we have

$$D_{\text{kl}}(P\|Q) \leq \log \int \frac{dP^2}{dQ} = \log(1 + D_{\chi^2}(P\|Q)).$$

The last inequality is immediate as $\log(1 + t) \leq t$ for all $t > -1$. \square

It is also possible to relate mutual information between distributions to f -divergences, and even to bound the mutual information above and below by the Hellinger distance for certain problems. In

this case, we consider the following situation: let $V \in \{0, 1\}$ uniformly at random, and conditional on $V = v$, draw $X \sim P_v$ for some distribution P_v on a space \mathcal{X} . Then we have that

$$I(X; V) = \frac{1}{2}D_{\text{kl}}(P_0 \| \bar{P}) + \frac{1}{2}D_{\text{kl}}(P_1 \| \bar{P})$$

where $\bar{P} = \frac{1}{2}P_0 + \frac{1}{2}P_1$. The divergence measure on the right side of the preceding identity is a special case of the *Jenson-Shannon divergence*, defined for $\lambda \in [0, 1]$ by

$$D_{\text{js}, \lambda}(P \| Q) := \lambda D_{\text{kl}}(P \| \lambda P + (1 - \lambda)Q) + D_{\text{kl}}(Q \| \lambda P + (1 - \lambda)Q), \quad (2.2.12)$$

which is a symmetrized and bounded variant of the typical KL-divergence (we use the shorthand $D_{\text{js}}(P \| Q) := D_{\text{js}, \frac{1}{2}}(P \| Q)$ for the symmetric case). As a consequence, we also have

$$I(X; V) = \frac{1}{2}D_f(P_0 \| P_1) + \frac{1}{2}D_f(P_1 \| P_0),$$

where $f(t) = -t \log(\frac{1}{2t} + \frac{1}{2}) = t \log \frac{2t}{t+1}$, so that the mutual information is a particular f -divergence. This form—as we see in the later chapters—is frequently convenient because it gives an object with similar tensorization properties to KL-divergence while enjoying the boundedness properties of Hellinger and variation distances. The following proposition captures the latter properties.

Proposition 2.2.10. *Let (X, V) be distributed as above. Then*

$$\log 2 \cdot d_{\text{hel}}^2(P_0, P_1) \leq I(X; V) = D_{\text{js}}(P_0 \| P_1) \leq \min \left\{ \frac{\log 2 \cdot \|P_0 - P_1\|_{\text{TV}}}{2 \cdot d_{\text{hel}}^2(P_0, P_1)} \right\}.$$

Proof The lower bound and upper bound involving the variation distance both follow from analytic bounds on the binary entropy functional $h_2(p) = -p \log p - (1-p) \log(1-p)$. By expanding the mutual information and letting p_0 and p_1 be densities of P_0 and P_1 with respect to some base measure μ , we have

$$\begin{aligned} 2I(X; V) &= 2D_{\text{js}}(P_0 \| P_1) = \int p_0 \log \frac{2p_0}{p_0 + p_1} d\mu + \int p_1 \log \frac{2p_1}{p_0 + p_1} d\mu \\ &= 2 \log 2 + \int (p_0 + p_1) \left[\frac{p_0}{p_1 + p_1} \log \frac{p_0}{p_0 + p_1} + \frac{p_1}{p_1 + p_1} \log \frac{p_1}{p_0 + p_1} \right] d\mu \\ &= 2 \log 2 - \int (p_0 + p_1) h_2 \left(\frac{p_0}{p_1 + p_0} \right) d\mu. \end{aligned}$$

We claim that

$$2 \log 2 \cdot \min\{p, 1 - p\} \leq h_2(p) \leq 2 \log 2 \cdot \sqrt{p(1 - p)}$$

for all $p \in [0, 1]$ (see Exercises 2.17 and 2.18). Then the upper and lower bounds on the information become nearly immediate.

For the variation-based upper bound on $I(X; V)$, we use the lower bound $h_2(p) \geq 2 \log 2 \cdot \min\{p, 1 - p\}$ to write

$$\begin{aligned} \frac{2}{\log 2} I(X; V) &\leq 2 - \int (p_0(x) + p_1(x)) \min \left\{ \frac{p_0(x)}{p_0(x) + p_1(x)}, \frac{p_1(x)}{p_0(x) + p_1(x)} \right\} d\mu(x) \\ &= 2 - 2 \int \min\{p_0(x), p_1(x)\} d\mu(x) \\ &= 2 \int (p_1(x) - \min\{p_0(x), p_1(x)\}) d\mu(x) = 2 \int_{p_1 > p_0} (p_1(x) - p_0(x)) d\mu(x). \end{aligned}$$

But of course the final integral is $\|P_1 - P_0\|_{\text{TV}}$, giving $I(X; V) \leq \log 2 \|P_0 - P_1\|_{\text{TV}}$. Conversely, for the lower bound on $D_{\text{js}}(P_0 \| P_1)$, we use the upper bound $h_2(p) \leq 2 \log 2 \cdot \sqrt{p(1-p)}$ to obtain

$$\begin{aligned} \frac{1}{\log 2} I(X; V) &\geq 1 - \int (p_0 + p_1) \sqrt{\frac{p_0}{p_1 + p_0} \left(1 - \frac{p_0}{p_1 + p_0}\right)} d\mu \\ &= 1 - \int \sqrt{p_0 p_1} d\mu = \frac{1}{2} \int (\sqrt{p_0} - \sqrt{p_1})^2 d\mu = d_{\text{hel}}^2(P_0, P_1) \end{aligned}$$

as desired.

The Hellinger-based upper bound is simpler: by Proposition 2.2.9, we have

$$\begin{aligned} D_{\text{js}}(P_0 \| P_1) &= \frac{1}{2} D_{\text{kl}}(P_0 \| (P_0 + P_1)/2) + \frac{1}{2} D_{\text{kl}}(P_1 \| (P_0 + P_1)/2) \\ &\leq \frac{1}{2} D_{\chi^2}(P_0 \| (P_0 + P_1)/2) + \frac{1}{2} D_{\chi^2}(P_1 \| (P_0 + P_1)/2) \\ &= \frac{1}{2} \int \frac{(p_0 - p_1)^2}{p_0 + p_1} d\mu = \frac{1}{2} \int \frac{(\sqrt{p_0} - \sqrt{p_1})^2 (\sqrt{p_0} + \sqrt{p_1})^2}{p_0 + p_1} d\mu. \end{aligned}$$

Now note that $(a + b)^2 \leq 2a^2 + 2b^2$ for any $a, b \in \mathbb{R}$, and so $(\sqrt{p_0} + \sqrt{p_1})^2 \leq 2(p_0 + p_1)$, and thus the final integral has bound $\int (\sqrt{p_0} - \sqrt{p_1})^2 d\mu = 2d_{\text{hel}}^2(P_0, P_1)$. \square

2.2.5 Convexity and data processing for divergence measures

f -divergences satisfy a number of very useful properties, which we use repeatedly throughout the lectures. As the KL-divergence is an f -divergence, it of course satisfies these conditions; however, we state them in fuller generality, treating the KL-divergence results as special cases and corollaries.

We begin by exhibiting the general data processing properties and convexity properties of f -divergences, each of which specializes to KL divergence. We leave the proof of each of these as exercises. First, we show that f -divergences are jointly convex in their arguments.

Proposition 2.2.11. *Let P_1, P_2, Q_1, Q_2 be distributions on a set \mathcal{X} and $f : \mathbb{R}_+ \rightarrow \mathbb{R}$ be convex. Then for any $\lambda \in [0, 1]$,*

$$D_f(\lambda P_1 + (1 - \lambda)P_2 \| \lambda Q_1 + (1 - \lambda)Q_2) \leq \lambda D_f(P_1 \| Q_1) + (1 - \lambda)D_f(P_2 \| Q_2).$$

The proof of this proposition we leave as Exercise 2.11, which we treat as a consequence of the more general “log-sum” like inequalities of Exercise 2.8. It is, however, an immediate consequence of the fully specified definition (2.2.5) of an f -divergence, because $\text{pers}(f)$ is jointly convex. As an immediate corollary, we see that the same result is true for KL-divergence as well.

Corollary 2.2.12. *The KL-divergence $D_{\text{kl}}(P \| Q)$ is jointly convex in its arguments P and Q .*

We can also provide more general data processing inequalities for f -divergences, paralleling those for the KL-divergence. In this case, we consider random variables X and Z on spaces \mathcal{X} and \mathcal{Z} , respectively, and a Markov transition kernel K giving the Markov chain $X \rightarrow Z$. That is, $K(\cdot | x)$ is a probability distribution on \mathcal{Z} for each $x \in \mathcal{X}$, and conditioned on $X = x$, Z has distribution $K(\cdot | x)$ so that $K(A | x) = \mathbb{P}(Z \in A | X = x)$. Certainly, this includes the situation

when $Z = \phi(X)$ for some function ϕ , and more generally when $Z = \phi(X, U)$ for a function ϕ and some additional randomness U . For a distribution P on X , we then define the marginals

$$K_P(A) := \int_{\mathcal{X}} K(A, x) dP(x).$$

We then have the following proposition.

Proposition 2.2.13. *Let P and Q be distributions on X and let K be any Markov kernel. Then*

$$D_f(K_P \| K_Q) \leq D_f(P \| Q).$$

See Exercise 2.10 for a proof.

As a corollary, we obtain the following data processing inequality for KL-divergences, where we abuse notation to write $D_{\text{kl}}(X \| Y) = D_{\text{kl}}(P \| Q)$ for random variables $X \sim P$ and $Y \sim Q$.

Corollary 2.2.14. *Let $X, Y \in \mathcal{X}$ be random variables, let $U \in \mathcal{U}$ be independent of X and Y , and let $\phi : \mathcal{X} \times \mathcal{U} \rightarrow \mathcal{Z}$ for some spaces $\mathcal{X}, \mathcal{U}, \mathcal{Z}$. Then*

$$D_{\text{kl}}(\phi(X, U) \| \phi(Y, U)) \leq D_{\text{kl}}(X \| Y).$$

Thus, further processing of random variables can only bring them “closer” in the space of distributions; downstream processing of signals cannot make them further apart as distributions.

2.3 First steps into optimal procedures: testing inequalities

As noted in the introduction, a central benefit of the information theoretic tools we explore is that they allow us to certify the optimality of procedures—that no other procedure could (substantially) improve upon the one at hand. The main tools for these certifications are often inequalities governing the best possible behavior of a variety of statistical tests. Roughly, we put ourselves in the following scenario: nature chooses one of a possible set of (say) k worlds, indexed by probability distributions P_1, P_2, \dots, P_k , and conditional on nature’s choice of the world—the distribution $P^* \in \{P_1, \dots, P_k\}$ chosen—we observe data X drawn from P^* . Intuitively, it will be difficult to decide which distribution P_i is the true P^* if all the distributions are similar—the divergence between the P_i is small, or the information between X and P^* is negligible—and easy if the distances between the distributions P_i are large. With this outline in mind, we present two inequalities, and first examples of their application, to make concrete these connections to the notions of information and divergence defined in this section.

2.3.1 Le Cam’s inequality and binary hypothesis testing

The simplest instantiation of the above setting is the case when there are only two possible distributions, P_1 and P_2 , and our goal is to make a decision on whether P_1 or P_2 is the distribution generating data we observe. Concretely, suppose that nature chooses one of the distributions P_1 or P_2 at random, and let $V \in \{1, 2\}$ index this choice. Conditional on $V = v$, we then observe a sample X drawn from P_v . Denoting by \mathbb{P} the joint distribution of V and X , we have for any test $\Psi : \mathcal{X} \rightarrow \{1, 2\}$ that the probability of error is then

$$\mathbb{P}(\Psi(X) \neq V) = \frac{1}{2}P_1(\Psi(X) \neq 1) + \frac{1}{2}P_2(\Psi(X) \neq 2).$$

We can give an exact expression for the minimal possible error in the above hypothesis test. Indeed, a standard result of Le Cam (see [134, 194, Lemma 1]) is the following variational representation of the total variation distance (2.2.6), which is the f -divergence associated with $f(t) = \frac{1}{2}|t-1|$, as a function of testing error.

Proposition 2.3.1. *Let \mathcal{X} be an arbitrary set. For any distributions P_1 and P_2 on \mathcal{X} , we have*

$$\inf_{\Psi} \{P_1(\Psi(X) \neq 1) + P_2(\Psi(X) \neq 2)\} = 1 - \|P_1 - P_2\|_{\text{TV}},$$

where the infimum is taken over all tests $\Psi : \mathcal{X} \rightarrow \{1, 2\}$.

Proof Any test $\Psi : \mathcal{X} \rightarrow \{1, 2\}$ has an acceptance region, call it $A \subset \mathcal{X}$, where it outputs 1 and a region A^c where it outputs 2.

$$P_1(\Psi \neq 1) + P_2(\Psi \neq 2) = P_1(A^c) + P_2(A) = 1 - P_1(A) + P_2(A).$$

Taking an infimum over such acceptance regions, we have

$$\inf_{\Psi} \{P_1(\Psi \neq 1) + P_2(\Psi \neq 2)\} = \inf_{A \subset \mathcal{X}} \{1 - (P_1(A) - P_2(A))\} = 1 - \sup_{A \subset \mathcal{X}} (P_1(A) - P_2(A)),$$

which yields the total variation distance as desired. \square

In the two-hypothesis case, we also know that the optimal test, by the Neyman-Pearson lemma, is a likelihood ratio test. That is, assuming that P_1 and P_2 have densities p_1 and p_2 , the optimal test is of the form

$$\Psi(X) = \begin{cases} 1 & \text{if } \frac{p_1(X)}{p_2(X)} \geq t \\ 2 & \text{if } \frac{p_1(X)}{p_2(X)} < t \end{cases}$$

for some threshold $t \geq 0$. In the case that the prior probabilities on P_1 and P_2 are each $\frac{1}{2}$, then $t = 1$ is optimal.

We give one example application of Proposition 2.3.1 to the problem of testing a normal mean.

Example 2.3.2 (Testing a normal mean): Suppose we observe $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P$ for $P = P_1$ or $P = P_2$, where P_v is the normal distribution $\mathcal{N}(\mu_v, \sigma^2)$, where $\mu_1 \neq \mu_2$. We would like to understand the sample size n necessary to guarantee that no test can have small error, that is, say, that

$$\inf_{\Psi} \{P_1(\Psi(X_1, \dots, X_n) \neq 1) + P_2(\Psi(X_1, \dots, X_n) \neq 2)\} \geq \frac{1}{2}.$$

By Proposition 2.3.1, we have that

$$\inf_{\Psi} \{P_1(\Psi(X_1, \dots, X_n) \neq 1) + P_2(\Psi(X_1, \dots, X_n) \neq 2)\} \geq 1 - \|P_1^n - P_2^n\|_{\text{TV}},$$

where P_v^n denotes the n -fold product of P_v , that is, the distribution of $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P_v$.

The interaction between total variation distance and product distributions is somewhat subtle, so it is often advisable to use a divergence measure more attuned to the i.i.d. nature of the sampling scheme. Two such measures are the KL-divergence and Hellinger distance, both of which we explore in the coming chapters. With that in mind, we apply Pinsker's

inequality (2.2.10) to see that $\|P_1^n - P_2^n\|_{\text{TV}}^2 \leq \frac{1}{2} D_{\text{kl}}(P_1^n \| P_2^n) = \frac{n}{2} D_{\text{kl}}(P_1 \| P_2)$, which implies that

$$1 - \|P_1^n - P_2^n\|_{\text{TV}} \geq 1 - \sqrt{\frac{n}{2} D_{\text{kl}}(P_1 \| P_2)} = 1 - \sqrt{\frac{n}{2} \left(\frac{1}{2\sigma^2} (\mu_1 - \mu_2)^2 \right)^{\frac{1}{2}}} = 1 - \frac{\sqrt{n} |\mu_1 - \mu_2|}{2\sigma}.$$

In particular, if $n \leq \frac{\sigma^2}{(\mu_1 - \mu_2)^2}$, then we have our desired lower bound of $\frac{1}{2}$.

Conversely, a calculation yields that $n \geq \frac{C\sigma^2}{(\mu_1 - \mu_2)^2}$, for some numerical constant $C \geq 1$, implies small probability of error. We leave this calculation to the reader. \diamond

2.3.2 Fano's inequality and multiple hypothesis testing

There are of course situations in which we do not wish to simply test two hypotheses, but have multiple hypotheses present. In such situations, Fano's inequality, which we present shortly, is the most common tool for proving fundamental limits, lower bounds on probability of error, and converses (to results on achievability of some performance level) in information theory. We write this section in terms of general random variables, ignoring the precise setting of selecting an index in a family of distributions, though that is implicit in what we do.

Let X be a random variable taking values in a finite set \mathcal{X} , and assume that we observe a (different) random variable Y , and then must estimate or guess the true value of \hat{X} . That is, we have the Markov chain

$$X \rightarrow Y \rightarrow \hat{X},$$

and we wish to provide lower bounds on the probability of error—that is, that $\hat{X} \neq X$. If we let the function $h_2(p) = -p \log p - (1-p) \log(1-p)$ denote the binary entropy (entropy of a Bernoulli random variable with parameter p), Fano's inequality takes the following form [e.g. 57, Chapter 2]:

Proposition 2.3.3 (Fano inequality). *For any Markov chain $X \rightarrow Y \rightarrow \hat{X}$, we have*

$$h_2(\mathbb{P}(\hat{X} \neq X)) + \mathbb{P}(\hat{X} \neq X) \log(|\mathcal{X}| - 1) \geq H(X | \hat{X}). \quad (2.3.1)$$

Proof This proof follows by expanding an entropy functional in two different ways. Let E be the indicator for the event that $\hat{X} \neq X$, that is, $E = 1$ if $\hat{X} \neq X$ and is 0 otherwise. Then we have

$$\begin{aligned} H(X, E | \hat{X}) &= H(X | E, \hat{X}) + H(E | \hat{X}) \\ &= \mathbb{P}(E = 1) H(X | E = 1, \hat{X}) + \mathbb{P}(E = 0) \underbrace{H(X | E = 0, \hat{X})}_{=0} + H(E | \hat{X}), \end{aligned}$$

where the zero follows because given there is no error, X has no variability given \hat{X} . Expanding the entropy by the chain rule in a different order, we have

$$H(X, E | \hat{X}) = H(X | \hat{X}) + \underbrace{H(E | \hat{X}, X)}_{=0},$$

because E is perfectly predicted by \hat{X} and X . Combining these equalities, we have

$$H(X | \hat{X}) = H(X, E | \hat{X}) = \mathbb{P}(E = 1) H(X | E = 1, \hat{X}) + H(E | \hat{X}).$$

Noting that $H(E | X) \leq H(E) = h_2(\mathbb{P}(E = 1))$, as conditioning reduces entropy, and that $H(X | E = 1, \hat{X}) \leq \log(|\mathcal{X}| - 1)$, as X can take on at most $|\mathcal{X}| - 1$ values when there is an error,

completes the proof. \square

We can rewrite Proposition 2.3.3 in a convenient way when X is uniform in \mathcal{X} . Indeed, by definition of the mutual information, we have $I(X; \hat{X}) = H(X) - H(X | \hat{X})$, so Proposition 9.4.1 implies that in the canonical hypothesis testing problem from Section 9.2.1, we have

Corollary 2.3.4. *Assume that X is uniform on \mathcal{X} . For any Markov chain $X \rightarrow Y \rightarrow \hat{X}$,*

$$\mathbb{P}(\hat{X} \neq X) \geq 1 - \frac{I(X; Y) + \log 2}{\log(|\mathcal{X}|)}. \quad (2.3.2)$$

Proof Let $P_{\text{error}} = \mathbb{P}(X \neq \hat{X})$ denote the probability of error. Noting that $h_2(p) \leq \log 2$ for any $p \in [0, 1]$ (recall inequality (2.1.2), that is, that uniform random variables maximize entropy), then using Proposition 9.4.1, we have

$$\log 2 + P_{\text{error}} \log(|\mathcal{X}|) \geq h_2(P_{\text{error}}) + P_{\text{error}} \log(|\mathcal{X}| - 1) \stackrel{(i)}{\geq} H(X | \hat{X}) \stackrel{(ii)}{=} H(X) - I(X; \hat{X}).$$

Here step (i) uses Proposition 2.3.3 and step (ii) uses the definition of mutual information, that $I(X; \hat{X}) = H(X) - H(X | \hat{X})$. The data processing inequality implies that $I(X; \hat{X}) \leq I(X; Y)$, and using $H(X) = \log(|\mathcal{X}|)$ completes the proof. \square

In particular, Corollary 2.3.4 shows that when X is chosen uniformly at random and we observe Y , we have

$$\inf_{\Psi} \mathbb{P}(\Psi(Y) \neq X) \geq 1 - \frac{I(X; Y) + \log 2}{\log |\mathcal{X}|},$$

where the infimum is taken over all testing procedures Ψ . Some interpretation of this quantity is helpful. If we think roughly of the number of bits it takes to describe a variable X uniformly chosen from \mathcal{X} , then we expect that $\log_2 |\mathcal{X}|$ bits are necessary (and sufficient). Thus, until we collect enough information that $I(X; Y) \approx \log |\mathcal{X}|$, so that $I(X; Y) / \log |\mathcal{X}| \approx 1$, we are unlikely to be unable to identify the variable X with any substantial probability. So we must collect enough bits to actually discover X .

Example 2.3.5 (20 questions game): In the 20 questions game—a standard children’s game—there are two players, the “chooser” and the “guesser,” and an agreed upon universe \mathcal{X} . The chooser picks an element $x \in \mathcal{X}$, and the guesser’s goal is to find x by using a series of yes/no questions about x . We consider optimal strategies for each player in this game, assuming that \mathcal{X} is finite and letting $m = |\mathcal{X}|$ be the universe size for shorthand.

For the guesser, it is clear that at most $\lceil \log_2 m \rceil$ questions are necessary to guess the item X that the chooser has picked—at each round of the game, the guesser asks a question that eliminates half of the remaining possible items. Indeed, let us assume that $m = 2^l$ for some $l \in \mathbb{N}$; if not, the guesser can always make her task more difficult by increasing the size of \mathcal{X} until it is a power of 2. Thus, after k rounds, there are $m2^{-k}$ items left, and we have

$$m \left(\frac{1}{2} \right)^k \leq 1 \quad \text{if and only if} \quad k \geq \log_2 m.$$

For the converse—the chooser’s strategy—let Y_1, Y_2, \dots, Y_k be the sequence of yes/no answers given to the guesser. Assume that the chooser picks X uniformly at random in \mathcal{X} . Then Fano’s inequality (2.3.2) implies that for the guess \hat{X} the guesser makes,

$$\mathbb{P}(\hat{X} \neq X) \geq 1 - \frac{I(X; Y_1, \dots, Y_k) + \log 2}{\log m}.$$

By the chain rule for mutual information, we have

$$I(X; Y_1, \dots, Y_k) = \sum_{i=1}^k I(X; Y_i \mid Y_{1:i-1}) = \sum_{i=1}^k H(Y_i \mid Y_{1:i-1}) - H(Y_i \mid Y_{1:i-1}, X) \leq \sum_{i=1}^k H(Y_i).$$

As the answers Y_i are yes/no, we have $H(Y_i) \leq \log 2$, so that $I(X; Y_{1:k}) \leq k \log 2$. Thus we find

$$\mathbb{P}(\hat{X} \neq X) \geq 1 - \frac{(k+1) \log 2}{\log m} = \frac{\log_2 m - 1}{\log_2 m} - \frac{k}{\log_2 m},$$

so that the guesser must have $k \geq \log_2(m/2)$ to be guaranteed that she will make no mistakes. \diamond

2.4 A first operational result: entropy and source coding

The final section of this chapter explores the basic results in source coding. Source coding—in its simplest form—tells us precisely the number of bits (or some other form of information storage) are necessary to perfectly encode a sequence of random variables X_1, X_2, \dots drawn according to a known distribution P .

2.4.1 The source coding problem

Assume we receive data consisting of a sequence of symbols X_1, X_2, \dots , drawn from a known distribution P on a finite or countable space \mathcal{X} . We wish to choose an encoding, represented by a *d-ary code function* C that maps \mathcal{X} to finite strings consisting of the symbols $\{0, 1, \dots, d-1\}$. We denote this by $C : \mathcal{X} \rightarrow \{0, 1, \dots, d-1\}^*$, where the superscript $*$ denotes the length may change from input to input, and use $\ell_C(x)$ to denote the length of the string $C(x)$.

In general, we will consider a variety of types of codes; we define each in order of complexity of their decoding.

Definition 2.1. A *d-ary code* $C : \mathcal{X} \rightarrow \{0, \dots, d-1\}^*$ is non-singular if for each $x, x' \in \mathcal{X}$ we have

$$C(x) \neq C(x') \quad \text{if } x \neq x'.$$

While Definition 2.1 is natural, generally speaking, we wish to transmit or encode a variety of code-words simultaneously, that is, we wish to encode a sequence X_1, X_2, \dots using the natural *extension* of the code C as the string $C(X_1)C(X_2)C(X_3)\dots$, where $C(x_1)C(x_2)$ denotes the concatenation of the strings $C(x_1)$ and $C(x_2)$. In this case, we require that the code be uniquely decodable:

Definition 2.2. A *d-ary code* $C : \mathcal{X} \rightarrow \{0, \dots, d-1\}^*$ is uniquely decodable if for all sequences $x_1, \dots, x_n \in \mathcal{X}$ and $x'_1, \dots, x'_n \in \mathcal{X}$ we have

$$C(x_1)C(x_2)\dots C(x_n) = C(x'_1)C(x'_2)\dots C(x'_n) \quad \text{if and only if } x_1 = x'_1, \dots, x_n = x'_n.$$

That is, the extension of the code C to sequences is non-singular.

While more useful (generally) than simply non-singular codes, uniquely decodable codes may require inspection of an entire string before recovering the first element. With that in mind, we now consider the easiest to use codes, which can always be decoded instantaneously.

Definition 2.3. A d -ary code $C : \mathcal{X} \rightarrow \{0, \dots, d-1\}^*$ is uniquely decodable or instantaneous if no codeword is the prefix to another codeword.

As is hopefully apparent from the definitions, all prefix/instantaneous codes are uniquely decodable, which are in turn non-singular. The converse is not true, though we will see a sense in which—as long as we care only about encoding sequences—using prefix instead of uniquely decodable codes has negligible consequences.

For example, written English, with periods (.) and spaces () included at the ends of words (among other punctuation) is an instantaneous encoding of English into the symbols of the alphabet and punctuation, as punctuation symbols enforce that no “codeword” is a prefix of any other. A few more concrete examples may make things more clear.

Example 2.4.1 (Encoding strategies): Consider the encoding schemes below, which encode the letters a, b, c, and d.

Symbol	$C_1(x)$	$C_2(x)$	$C_3(x)$
a	0	00	0
b	00	10	10
c	000	11	110
d	0000	110	111

By inspection, it is clear that C_1 is non-singular but certainly not uniquely decodable (does the sequence 0000 correspond to aaaa, bb, aab, aba, baa, ca, ac, or d?), while C_3 is a prefix code. We leave showing that C_2 is uniquely decodable as an exercise. \diamond

2.4.2 The Kraft-McMillan inequalities

We now turn to a few results on the connections between source-coding and entropy. Our first result, the *Kraft-McMillan inequality*, is an essential result that—as we shall see—essentially says that there is no difference in code-lengths attainable by prefix codes and uniquely decodable codes.

Theorem 2.4.2. Let \mathcal{X} be a finite or countable set, and let $\ell : \mathcal{X} \rightarrow \mathbb{N}$ be a function. If $\ell(x)$ is the length of the encoding of the symbol x in a uniquely decodable d -ary code, then

$$\sum_{x \in \mathcal{X}} d^{-\ell(x)} \leq 1. \quad (2.4.1)$$

Conversely, given any function $\ell : \mathcal{X} \rightarrow \mathbb{N}$ satisfying inequality (2.4.1), there is a prefix code whose codewords have length $\ell(x)$ for each $x \in \mathcal{X}$.

Proof We prove the first statement of the theorem first by a counting and asymptotic argument.

We begin by assuming that \mathcal{X} is finite; we eliminate this assumption subsequently. As a consequence, there is some maximum length ℓ_{\max} such that $\ell(x) \leq \ell_{\max}$ for all $x \in \mathcal{X}$. For a sequence $x_1, \dots, x_n \in \mathcal{X}$, we have by the definition of our encoding strategy that $\ell(x_1, \dots, x_n) = \sum_{i=1}^n \ell(x_i)$. In addition, for each m we let

$$E_n(m) := \{x_{1:n} \in \mathcal{X}^n \text{ such that } \ell(x_{1:n}) = m\}$$

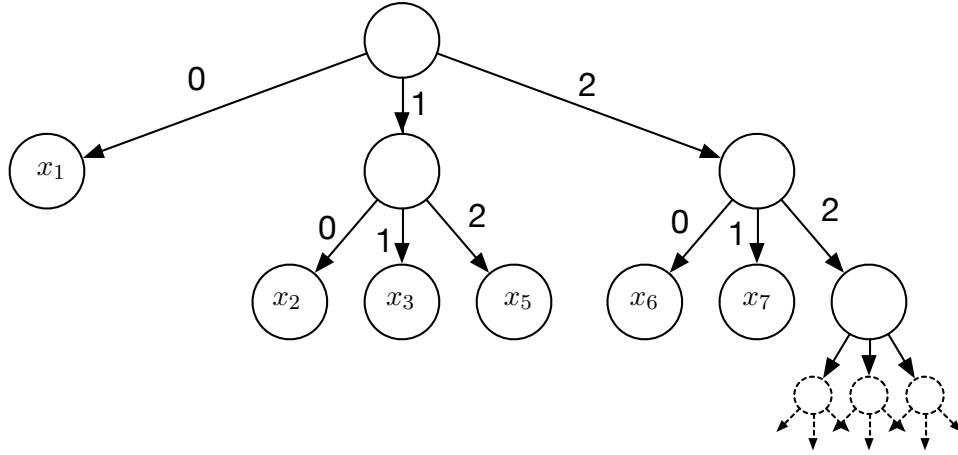


Figure 2.1. Prefix-tree encoding of a set of symbols. The encoding for x_1 is 0, for x_2 is 10, for x_3 is 11, for x_4 is 12, for x_5 is 20, for x_6 is 21, and nothing is encoded as 1, 2, or 22.

denote the symbols x encoded with codewords of length m in our code, then as the code is uniquely decodable we certainly have $\text{card}(E_n(m)) \leq d^m$ for all n and m . Moreover, for all $x_{1:n} \in \mathcal{X}^n$ we have $\ell(x_{1:n}) \leq n\ell_{\max}$. We thus re-index the sum $\sum_x d^{-\ell(x)}$ and compute

$$\begin{aligned} \sum_{x_1, \dots, x_n \in \mathcal{X}^n} d^{-\ell(x_1, \dots, x_n)} &= \sum_{m=1}^{n\ell_{\max}} \text{card}(E_n(m)) d^{-m} \\ &\leq \sum_{m=1}^{n\ell_{\max}} d^{m-m} = n\ell_{\max}. \end{aligned}$$

The preceding relation is true for all $n \in \mathbb{N}$, so that

$$\left(\sum_{x_{1:n} \in \mathcal{X}^n} d^{-\ell(x_{1:n})} \right)^{1/n} \leq n^{1/n} \ell_{\max}^{1/n} \rightarrow 1$$

as $n \rightarrow \infty$. In particular, using that

$$\sum_{x_{1:n} \in \mathcal{X}^n} d^{-\ell(x_{1:n})} = \sum_{x_1, \dots, x_n \in \mathcal{X}^n} d^{-\ell(x_1)} \dots d^{-\ell(x_n)} = \left(\sum_{x \in \mathcal{X}} d^{-\ell(x)} \right)^n,$$

we obtain $\sum_{x \in \mathcal{X}} d^{-\ell(x)} \leq 1$.

Returning to the case that $\text{card}(\mathcal{X}) = \infty$, by defining the sequence

$$D_k := \sum_{x \in \mathcal{X}, \ell(x) \leq k} d^{-\ell(x)},$$

as each subset $\{x \in \mathcal{X} : \ell(x) \leq k\}$ is uniquely decodable, we have $D_k \leq 1$ for all k . Then $1 \geq \lim_{k \rightarrow \infty} D_k = \sum_{x \in \mathcal{X}} d^{-\ell(x)}$.

The achievability of such a code intuitively follows by a pictorial argument (recall Figure 2.1), so we first sketch the result non-rigorously. Indeed, let \mathcal{T}_d be an (infinite) d -ary tree. Then, at each

level m of the tree, assign one of the nodes at that level to each symbol $x \in \mathcal{X}$ such that $\ell(x) = m$. Eliminate the subtree below that node, and repeat with the remaining symbols. The codeword corresponding to symbol x is then the path to the symbol in the tree.

A more formal version implementing this sketch follows. Let ℓ be a length function satisfying $\sum_{x \in \mathcal{X}} d^{-\ell(x)} \leq 1$. Identify \mathcal{X} with \mathbb{N} (or a subset thereof) in such a way that $1 \leq \ell(1) \leq \ell(2) \leq \dots$, i.e., $\ell(x) \leq \ell(y)$ whenever $x < y$, and let $\mathcal{X}_m = \{x \in \mathcal{X} \mid \ell(x) = m\}$ be the set of inputs with encoding length m . For each $x \in \mathbb{N}$, define the value

$$v(x) = \sum_{i < x} d^{-\ell(i)}.$$

We let the codeword $C(x)$ for x be the first $\ell(x)$ terms in the d -ary expansion of $v(x)$. Certainly the length of this encoding satisfies $|C(x)| = \ell(x)$. To see that it is prefix-free, take two symbols $x < y$, and assume for the sake of contradiction that $C(x)$ is a prefix of $C(y)$. Then $v(y) \geq v(x)$, while $v(y) - v(x) \leq d^{-\ell(x)}$ because the two representations agree on the first $\ell(x)$ terms in the expansion. But

$$v(y) - v(x) = \sum_{i < y} d^{-\ell(i)} - \sum_{i < x} d^{-\ell(i)} = \sum_{x \leq i < y} d^{-\ell(i)} = d^{-\ell(x)} + \sum_{x < i < y} d^{-\ell(i)} > d^{-\ell(x)},$$

a contradiction. □

With the Kraft-McMillan theorem in place, we may directly relate the entropy of a random variable to the length of possible encodings for the variable; in particular, we show that the entropy is essentially *the best* possible code length of a uniquely decodable source code. In this theorem, we use the shorthand

$$H_d(X) := - \sum_{x \in \mathcal{X}} p(x) \log_d p(x).$$

Theorem 2.4.3. *Let $X \in \mathcal{X}$ be a discrete random variable distributed according to P and let ℓ_C be the length function associated with a d -ary encoding $C : \mathcal{X} \rightarrow \{0, \dots, d-1\}^*$. In addition, let \mathcal{C} be the set of all uniquely decodable d -ary codes for \mathcal{X} . Then*

$$H_d(X) \leq \inf \{ \mathbb{E}_P[\ell_C(X)] : C \in \mathcal{C} \} \leq H_d(X) + 1.$$

Proof The lower bound is an argument by convex optimization, while for the upper bound we give an explicit length function and (implicit) prefix code attaining the bound. For the lower bound, we assume for simplicity that \mathcal{X} is finite, and we identify $\mathcal{X} = \{1, \dots, |\mathcal{X}|\}$ (let $m = |\mathcal{X}|$ for shorthand). Then as \mathcal{C} consists of *uniquely decodable* codebooks, all the associated length functions must satisfy the Kraft-McMillan inequality (2.4.1). Letting $\ell_i = \ell(i)$, the minimal encoding length is at least

$$\inf_{\ell \in \mathbb{R}^m} \left\{ \sum_{i=1}^m p_i \ell_i : \sum_{i=1}^m d^{-\ell_i} \leq 1 \right\}.$$

By introducing the Lagrange multiplier $\lambda \geq 0$ for the inequality constraint, we may write the Lagrangian for the preceding minimization problem as

$$\mathcal{L}(\ell, \lambda) = p^\top \ell + \lambda \left(\sum_{i=1}^n d^{-\ell_i} - 1 \right) \quad \text{with} \quad \nabla_\ell \mathcal{L}(\ell, \lambda) = p - \lambda \left[d^{-\ell_i} \log d \right]_{i=1}^m.$$

In particular, the optimal ℓ satisfies $\ell_i = \log_d \frac{\theta}{p_i}$ for some constant θ , and solving $\sum_{i=1}^m d^{-\log_d \frac{\theta}{p_i}} = 1$ gives $\theta = 1$ and $\ell(i) = \log_d \frac{1}{p_i}$.

To attain the result, simply set our encoding to be $\ell(x) = \left\lceil \log_d \frac{1}{P(X=x)} \right\rceil$, which satisfies the Kraft-McMillan inequality and thus yields a valid prefix code with

$$\mathbb{E}_P[\ell(X)] = \sum_{x \in \mathcal{X}} p(x) \left\lceil \log_d \frac{1}{p(x)} \right\rceil \leq - \sum_{x \in \mathcal{X}} p(x) \log_d p(x) + 1 = H_d(X) + 1$$

as desired. \square

Theorem 2.4.3 thus shows that, at least to within an additive constant of 1, the entropy both upper and lower bounds the expected length of a uniquely decodable code for the random variable X . This is the first of our promised “operational interpretations” of the entropy.

2.4.3 Entropy rates and longer codes

Theorem 2.4.3 is a bit unsatisfying in that the additive constant 1 may be quite large relative to the entropy. By allowing encoding longer sequences, we can (asymptotically) eliminate this error factor. To that end, we here show that it is possible, at least for appropriate distributions on random variables X_i , to achieve a per-symbol encoding length that approaches a limiting version of the Shannon entropy of a random variable. We give two definitions capturing the limiting entropy properties of sequences of random variables.

Definition 2.4. *The entropy rate of a sequence X_1, X_2, \dots of random variables is*

$$H(\{X_i\}) := \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n) \quad (2.4.2)$$

whenever the limit exists.

In some situations, the limit (2.4.2) may not exist. However, there are a variety of situations in which it does, and we focus generally on a specific but common instance in which the limit does exist. First, we recall the definition of a stationary sequence of random variables.

Definition 2.5. *We say a sequence X_1, X_2, \dots of random variable is stationary if for all n and all $k \in \mathbb{N}$ and all measurable sets $A_1, \dots, A_k \subset \mathcal{X}$ we have*

$$\mathbb{P}(X_1 \in A_1, \dots, X_k \in A_k) = \mathbb{P}(X_{n+1} \in A_1, \dots, X_{n+k} \in A_k).$$

With this definition, we have the following result.

Proposition 2.4.4. *Let the sequence of random variables $\{X_i\}$, taking values in the discrete space \mathcal{X} , be stationary. Then*

$$H(\{X_i\}) = \lim_{n \rightarrow \infty} H(X_n \mid X_1, \dots, X_{n-1})$$

and the limits (2.4.2) and above exist.

Proof We begin by making the following standard observation of Cesàro means: if $c_n = \frac{1}{n} \sum_{i=1}^n a_i$ and $a_i \rightarrow a$, then $c_n \rightarrow a$.³ Now, we note that for a stationary sequence, we have that

$$H(X_n | X_{1:n-1}) = H(X_{n+1} | X_{2:n}),$$

and using that conditioning decreases entropy, we have

$$H(X_{n+1} | X_{1:n}) \leq H(X_n | X_{1:n-1}).$$

Thus the sequence $a_n := H(X_n | X_{1:n-1})$ is non-increasing and bounded below by 0, so that it has some limit $\lim_{n \rightarrow \infty} H(X_n | X_{1:n-1})$. As $H(X_1, \dots, X_n) = \sum_{i=1}^n H(X_i | X_{1:i-1})$ by the chain rule for entropy, we achieve the result of the proposition. \square

Finally, we present a result showing that it is possible to achieve average code length of at most the entropy rate, which for stationary sequences is smaller than the entropy of any single random variable X_i . To do so, we require the use of a block code, which (while it may be prefix code) treats sets of random variables $(X_1, \dots, X_m) \in \mathcal{X}^m$ as a single symbol to be jointly encoded.

Proposition 2.4.5. *Let the sequence of random variables X_1, X_2, \dots be stationary. Then for any $\epsilon > 0$, there exists an $m \in \mathbb{N}$ and a d -ary (prefix) block encoder $C : \mathcal{X}^m \rightarrow \{0, \dots, d-1\}^*$ such that*

$$\lim_n \frac{1}{n} \mathbb{E}_P[\ell_C(X_{1:n})] \leq H(\{X_i\}) + \epsilon = \lim_n H(X_n | X_1, \dots, X_{n-1}) + \epsilon.$$

Proof Let $C : \mathcal{X}^m \rightarrow \{0, 1, \dots, d-1\}^*$ be any prefix code with

$$\ell_C(x_{1:m}) \leq \left\lceil \log \frac{1}{P(X_{1:m} = x_{1:m})} \right\rceil.$$

Then whenever n/m is an integer, we have

$$\begin{aligned} \mathbb{E}_P[\ell_C(X_{1:n})] &= \sum_{i=1}^{n/m} \mathbb{E}_P[\ell_C(X_{mi+1}, \dots, X_{m(i+1)})] \leq \sum_{i=1}^{n/m} [H(X_{mi+1}, \dots, X_{m(i+1)}) + 1] \\ &= \frac{n}{m} + \frac{n}{m} H(X_1, \dots, X_m). \end{aligned}$$

Dividing by n gives the result by taking m suitably large that $\frac{1}{m} + \frac{1}{m} H(X_1, \dots, X_m) \leq \epsilon + H(\{X_i\})$.

Note that if the m does not divide n , we may also encode the length of the sequence of encoded words in each block of length m ; in particular, if the block begins with a 0, it encodes m symbols, while if it begins with a 1, then the next $\lceil \log_d m \rceil$ bits encode the length of the block. This would yield an increase in the expected length of the code to

$$\mathbb{E}_P[\ell_C(X_{1:n})] \leq \frac{2n + \lceil \log_2 m \rceil}{m} + \frac{n}{m} H(X_1, \dots, X_m).$$

Dividing by n and letting $n \rightarrow \infty$ gives the result, as we can always choose m large. \square

³Indeed, let $\epsilon > 0$ and take N such that $n \geq N$ implies that $|a_i - a| < \epsilon$. Then for $n \geq N$, we have

$$c_n - a = \frac{1}{n} \sum_{i=1}^n (a_i - a) = \frac{N(c_N - a)}{n} + \frac{1}{n} \sum_{i=N+1}^n (a_i - a) \in \frac{N(c_N - a)}{n} \pm \epsilon.$$

Taking $n \rightarrow \infty$ yields that the term $N(c_N - a)/n \rightarrow 0$, which gives that $c_n - a \in [-\epsilon, \epsilon]$ eventually for any $\epsilon > 0$, which is our desired result.

2.5 Bibliography

The material in this chapter is classical in information theory. For all of our treatment of mutual information, entropy, and KL-divergence in the discrete case, [Cover and Thomas](#) provide an essentially complete treatment in Chapter 2 of their book [57]. Gray [104] provides a more advanced (measure-theoretic) version of these results, with Chapter 5 covering most of our results (or Chapter 7 in the newer addition of the same book). Csiszár and Körner [61] is the classic reference for coding theorems and results on communication, including stronger converse results.

The f -divergence was independently discovered by Ali and Silvey [6] and Csiszár [59], and is consequently sometimes called an Ali-Silvey divergence or Csiszár divergence. Liese and Vajda [137] provide a survey of f -divergences and their relationships with different statistical concepts (taking a Bayesian point of view), and various authors have extended the pairwise divergence measures to divergence measures between multiple distributions [107], making connections to experimental design and classification [98, 76], which we investigate later in book. The inequalities relating divergences in Section 2.2.4 are now classical, and standard references present them [134, 182]. For a proof that equality (2.2.4) is equivalent to the definition (2.2.3) with the appropriate closure operations, see the paper [76, Proposition 1]. We borrow the proof of the upper bound in Proposition 2.2.10 from the paper [138].

JCD Comment: Converse to Kraft is Chaitin?

2.6 Exercises

Our first few questions investigate properties of a divergence between distributions that is weaker than the KL-divergence, but is intimately related to optimal testing. Let P_1 and P_2 be arbitrary distributions on a space \mathcal{X} . The *total variation distance* between P_1 and P_2 is defined as

$$\|P_1 - P_2\|_{\text{TV}} := \sup_{A \subset \mathcal{X}} |P_1(A) - P_2(A)|.$$

Exercise 2.1: Prove the following identities about total variation. Throughout, let P_1 and P_2 have densities p_1 and p_2 on a (common) set \mathcal{X} .

- (a) $2 \|P_1 - P_2\|_{\text{TV}} = \int |p_1(x) - p_2(x)| dx.$
- (b) For functions $f : \mathcal{X} \rightarrow \mathbb{R}$, define the supremum norm $\|f\|_{\infty} = \sup_{x \in \mathcal{X}} |f(x)|$. Show that $2 \|P_1 - P_2\|_{\text{TV}} = \sup_{\|f\|_{\infty} \leq 1} \int_{\mathcal{X}} f(x)(p_1(x) - p_2(x)) dx.$
- (c) $\|P_1 - P_2\|_{\text{TV}} = \int \max\{p_1(x), p_2(x)\} dx - 1.$
- (d) $\|P_1 - P_2\|_{\text{TV}} = 1 - \int \min\{p_1(x), p_2(x)\} dx.$
- (e) For functions $f, g : \mathcal{X} \rightarrow \mathbb{R}$,

$$\inf \left\{ \int f(x)p_1(x)dx + \int g(x)p_2(x)dx : f + g \geq 1, f \geq 0, g \geq 0 \right\} = 1 - \|P_1 - P_2\|_{\text{TV}}.$$

Exercise 2.2 (Divergence between multivariate normal distributions): Let P_1 be $\mathbf{N}(\theta_1, \Sigma)$ and P_2 be $\mathbf{N}(\theta_2, \Sigma)$, where $\Sigma \succ 0$ is a positive definite matrix.

(a) Give $D_{\text{kl}}(P_1 \| P_2)$.

(b) Show that $d_{\text{hel}}^2(P_1, P_2) = 1 - \exp(-\frac{1}{8}(\mu_0 - \mu_1)^\top \Sigma^{-1}(\mu_0 - \mu_1))$.

Exercise 2.3 (The optimal test between distributions): Prove Le-Cam's inequality: for any function ψ with $\text{dom } \psi \supset \mathcal{X}$ and any distributions P_1, P_2 ,

$$P_1(\psi(X) \neq 1) + P_2(\psi(X) \neq 2) \geq 1 - \|P_1 - P_2\|_{\text{TV}}.$$

Thus, the sum of the probabilities of error in a hypothesis testing problem, where based on a sample X we must decide whether P_1 or P_2 is more likely, has value at least $1 - \|P_1 - P_2\|_{\text{TV}}$. Given P_1 and P_2 is this risk attainable?

Exercise 2.4: A random variable X has $\text{Laplace}(\lambda, \mu)$ distribution if it has density $p(x) = \frac{\lambda}{2} \exp(-\lambda|x-\mu|)$. Consider the hypothesis test of P_1 versus P_2 , where X has distribution $\text{Laplace}(\lambda, \mu_1)$ under P_1 and distribution $\text{Laplace}(\lambda, \mu_2)$ under P_2 , where $\mu_1 < \mu_2$. Show that the minimal value over all tests ψ of P_1 versus P_2 is

$$\inf_{\psi} \{P_1(\psi(X) \neq 1) + P_2(\psi(X) \neq 2)\} = \exp\left(-\frac{\lambda}{2}|\mu_1 - \mu_2|\right).$$

Exercise 2.5 (Log-sum inequality): Let a_1, \dots, a_n and b_1, \dots, b_n be non-negative reals. Show that

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \left(\sum_{i=1}^n a_i\right) \log \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i}.$$

(Hint: use the convexity of the function $x \mapsto -\log(x)$.)

Exercise 2.6: Given quantizers g_1 and g_2 , we say that g_1 is a *finer* quantizer than g_2 under the following condition: assume that g_1 induces the partition A_1, \dots, A_n and g_2 induces the partition B_1, \dots, B_m ; then for any of the sets B_i , there are exists some k and sets A_{i_1}, \dots, A_{i_k} such that $B_i = \cup_{j=1}^k A_{i_j}$. We let $g_1 \prec g_2$ denote that g_1 is a finer quantizer than g_2 . Prove

(a) Finer partitions increase the KL divergence: if $g_1 \prec g_2$,

$$D_{\text{kl}}(P \| Q | g_2) \leq D_{\text{kl}}(P \| Q | g_1).$$

(b) If \mathcal{X} is discrete (so P and Q have p.m.f.s p and q) then

$$D_{\text{kl}}(P \| Q) = \sum_x p(x) \log \frac{p(x)}{q(x)}.$$

Exercise 2.7 (f -divergences generalize standard divergences): Show the following properties of f -divergences:

(a) If $f(t) = |t - 1|$, then $D_f(P \| Q) = 2 \|P - Q\|_{\text{TV}}$.

(b) If $f(t) = t \log t$, then $D_f(P \| Q) = D_{\text{kl}}(P \| Q)$.

(c) If $f(t) = t \log t - \log t$, then $D_f(P \| Q) = D_{\text{kl}}(P \| Q) + D_{\text{kl}}(Q \| P)$.

(d) For any convex f satisfying $f(1) = 0$, $D_f(P\|Q) \geq 0$. (Hint: use Jensen's inequality.)

Exercise 2.8 (Generalized “log-sum” inequalities): Let $f : \mathbb{R}_+ \rightarrow \mathbb{R}$ be an arbitrary convex function.

(a) Let $a_i, b_i, i = 1, \dots, n$ be non-negative reals. Prove that

$$\left(\sum_{i=1}^n a_i \right) f \left(\frac{\sum_{i=1}^n b_i}{\sum_{i=1}^n a_i} \right) \leq \sum_{i=1}^n a_i f \left(\frac{b_i}{a_i} \right).$$

(b) Generalizing the preceding result, let $a : \mathcal{X} \rightarrow \mathbb{R}_+$ and $b : \mathcal{X} \rightarrow \mathbb{R}_+$, and let μ be a finite measure on \mathcal{X} with respect to which a is integrable. Show that

$$\int a(x) d\mu(x) f \left(\frac{\int b(x) d\mu(x)}{\int a(x) d\mu(x)} \right) \leq \int a(x) f \left(\frac{b(x)}{a(x)} \right) d\mu(x).$$

If you are unfamiliar with measure theory, prove the following essentially equivalent result: let $u : \mathcal{X} \rightarrow \mathbb{R}_+$ satisfy $\int u(x) dx < \infty$. Show that

$$\int a(x) u(x) dx f \left(\frac{\int b(x) u(x) dx}{\int a(x) u(x) dx} \right) \leq \int a(x) f \left(\frac{b(x)}{a(x)} \right) u(x) dx$$

whenever $\int a(x) u(x) dx < \infty$. (It is possible to demonstrate this remains true under appropriate limits even when $\int a(x) u(x) dx = +\infty$, but it is a mess.)

(Hint: use the fact that the perspective of a function f , defined by $h(x, t) = tf(x/t)$ for $t > 0$, is jointly convex in x and t (see Proposition B.3.12).

Exercise 2.9 (Data processing and f -divergences I): As with the KL-divergence, given a quantizer g of the set \mathcal{X} , where g induces a partition A_1, \dots, A_m of \mathcal{X} , we define the f -divergence between P and Q conditioned on g as

$$D_f(P\|Q \mid g) := \sum_{i=1}^m Q(A_i) f \left(\frac{P(A_i)}{Q(A_i)} \right) = \sum_{i=1}^m Q(g^{-1}(\{i\})) f \left(\frac{P(g^{-1}(\{i\}))}{Q(g^{-1}(\{i\}))} \right).$$

Given quantizers g_1 and g_2 , we say that g_1 is a *finer* quantizer than g_2 under the following condition: assume that g_1 induces the partition A_1, \dots, A_n and g_2 induces the partition B_1, \dots, B_m ; then for any of the sets B_i , there are exists some k and sets A_{i_1}, \dots, A_{i_k} such that $B_i = \cup_{j=1}^k A_{i_j}$. We let $g_1 \prec g_2$ denote that g_1 is a finer quantizer than g_2 .

(a) Let g_1 and g_2 be quantizers of the set \mathcal{X} , and let $g_1 \prec g_2$, meaning that g_1 is a finer quantization than g_2 . Prove that

$$D_f(P\|Q \mid g_2) \leq D_f(P\|Q \mid g_1).$$

Equivalently, show that whenever \mathcal{A} and \mathcal{B} are collections of sets partitioning \mathcal{X} , but \mathcal{A} is a finer partition of \mathcal{X} than \mathcal{B} , that

$$\sum_{B \in \mathcal{B}} Q(B) f \left(\frac{P(B)}{Q(B)} \right) \leq \sum_{A \in \mathcal{A}} Q(A) f \left(\frac{P(A)}{Q(A)} \right).$$

(Hint: Use the result of Question 2.8(a)).

- (b) Suppose that \mathcal{X} is countable (or finite) so that P and Q have p.m.f.s p and q . Show that

$$D_f(P\|Q) = \sum_x q(x) f\left(\frac{p(x)}{q(x)}\right),$$

where on the left we are using the partition definition (2.2.3); you should show that the partition into discrete parts of \mathcal{X} achieves the supremum. You may assume that \mathcal{X} is finite. (Though feel free to prove the result in the case that \mathcal{X} is infinite.)

Exercise 2.10 (General data processing inequalities): Let f be a convex function satisfying $f(1) = 0$. Let K be a Markov transition kernel from \mathcal{X} to \mathcal{Z} , that is, $K(\cdot, x)$ is a probability distribution on \mathcal{Z} for each $x \in \mathcal{X}$. (Written differently, we have $X \rightarrow Z$, and conditioned on $X = x$, Z has distribution $K(\cdot, x)$, so that $K(A, x)$ is the probability that $Z \in A$ given $X = x$.)

- (a) Define the marginals $K_P(A) = \int K(A, x)p(x)dx$ and $K_Q(A) = \int K(A, x)q(x)dx$. Show that

$$D_f(K_P\|K_Q) \leq D_f(P\|Q).$$

Hint: by equation (2.2.3), w.l.o.g. we may assume that \mathcal{Z} is finite and $\mathcal{Z} = \{1, \dots, m\}$; also recall Question 2.8.

- (b) Let X and Y be random variables with joint distribution P_{XY} and marginals P_X and P_Y . Define the f -information between X and Y as

$$I_f(X; Y) := D_f(P_{XY}\|P_X \times P_Y).$$

Use part (a) to show the following general data processing inequality: if we have the Markov chain $X \rightarrow Y \rightarrow Z$, then

$$I_f(X; Z) \leq I_f(X; Y).$$

Exercise 2.11 (Convexity of f -divergences): Prove Proposition 2.2.11. *Hint:* Use Question 2.8.

Exercise 2.12 (Variational forms of KL divergence): Let P and Q be arbitrary distributions on a common space \mathcal{X} . Prove the following variational representation, known as the Donsker-Varadhan theorem, of the KL divergence:

$$D_{\text{kl}}(P\|Q) = \sup_{f: \mathbb{E}_Q[e^{f(X)}] < \infty} \{\mathbb{E}_P[f(X)] - \log \mathbb{E}_Q[\exp(f(X))]\}.$$

You may assume that P and Q have densities.

Exercise 2.13: Let P and Q have densities p and q with respect to the base measure μ over the set \mathcal{X} . (Recall that this is no loss of generality, as we may take $\mu = P + Q$.) Define the support $\text{supp } P := \{x \in \mathcal{X} : p(x) > 0\}$. Show that

$$D_{\text{kl}}(P\|Q) \geq \log \frac{1}{Q(\text{supp } P)}.$$

Exercise 2.14: Let P_1 be $\mathcal{N}(\theta_1, \Sigma_1)$ and P_2 be $\mathcal{N}(\theta_2, \Sigma_2)$, where $\Sigma_i \succ 0$ are positive definite matrices. Give $D_{\text{kl}}(P_1\|P_2)$.

Exercise 2.15: Let $\{P_v\}_{v \in \mathcal{V}}$ be an arbitrary collection of distributions on a space \mathcal{X} and μ be a probability measure on \mathcal{V} . Show that if $V \sim \mu$ and conditional on $V = v$, we draw $X \sim P_v$, then

- (a) $I(X; V) = \int D_{\text{kl}}(P_v \| \bar{P}) d\mu(v)$, where $\bar{P} = \int P_v d\mu(v)$ is the (weighted) average of the P_v . You may assume that \mathcal{V} is discrete if you like.
- (b) For any distribution Q on \mathcal{X} , $I(X; V) = \int D_{\text{kl}}(P_v \| Q) d\mu(v) - D_{\text{kl}}(\bar{P} \| Q)$. Conclude that $I(X; V) \leq \int D_{\text{kl}}(P_v \| Q) d\mu(v)$, or, equivalently, \bar{P} minimizes $\int D_{\text{kl}}(P_v \| Q) d\mu(v)$ over all probabilities Q .

Exercise 2.16 (The triangle inequality for variation distance): Let P and Q be distributions on $X_1^n = (X_1, \dots, X_n) \in \mathcal{X}^n$, and let $P_i(\cdot | x_1^{i-1})$ be the conditional distribution of X_i given $X_1^{i-1} = x_1^{i-1}$ (and similarly for Q_i). Show that

$$\|P - Q\|_{\text{TV}} \leq \sum_{i=1}^n \mathbb{E}_P \left[\|P_i(\cdot | X_1^{i-1}) - Q_i(\cdot | X_1^{i-1})\|_{\text{TV}} \right],$$

where the expectation is taken over X_1^{i-1} distributed according to P .

Exercise 2.17: Let $h(p) = -p \log p - (1-p) \log(1-p)$. Show that $h(p) \geq 2 \log 2 \cdot \min\{p, 1-p\}$.

Exercise 2.18 (Lin [138], Theorem 8): Let $h(p) = -p \log p - (1-p) \log(1-p)$. Show that $h(p) \leq 2 \log 2 \cdot \sqrt{p(1-p)}$.

Exercise 2.19 (Proving Pinsker's inequality via data processing): We work through a proof of Proposition 2.2.8.(a) using the data processing inequality for f -divergences (Proposition 2.2.13).

- (a) Define $D_{\text{kl}}(p \| q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$. Argue that to prove Pinsker's inequality (2.2.10), it is enough to show that $(p-q)^2 \leq \frac{1}{2} D_{\text{kl}}(p \| q)$.
- (b) Define the negative binary entropy $h(p) = p \log p + (1-p) \log(1-p)$. Show that

$$h(p) \geq h(q) + h'(q)(p-q) + 2(p-q)^2$$

for any $p, q \in [0, 1]$.

- (c) Conclude Pinsker's inequality (2.2.10).

JCD Comment: Below are a few potential questions

Exercise 2.20: Use the paper “A New Metric for Probability Distributions” by Dominik Endres and Johannes Schindelin to prove that if $V \sim \text{Uniform}\{0, 1\}$ and $X | V = v \sim P_v$, then $\sqrt{I(X; V)}$ is a metric on distributions. (Said differently, $D_{\text{js}}(P \| Q)^{1/2}$ is a metric on distributions, and it generates the same topology as the TV-distance.)

Exercise 2.21: Relate the generalized Jensen-Shannon divergence between m distributions to redundancy in encoding.

Chapter 3

Exponential families and statistical modeling

Our second introductory chapter focuses on readers who may be less familiar with statistical modeling methodology and the how and why of fitting different statistical models. As in the preceding introductory chapter on information theory, this chapter will be a fairly terse blitz through the main ideas. Nonetheless, the ideas and distributions here should give us something on which to hang our hats, so to speak, as the distributions and models provide the basis for examples throughout the book. Exponential family models form the basis of much of statistics, as they are a natural step away from the most basic families of distributions—Gaussians—which admit exact computations but are brittle, to a more flexible set of models that retain enough analytical elegance to permit careful analyses while giving power in modeling. A key property is that fitting exponential family models reduces to the minimization of convex functions—convex optimization problems—an operation we treat as a technology akin to evaluating a function like \sin or \cos . This perspective (which is accurate enough) will arise throughout this book, and informs the philosophy we adopt that once we formulate a problem as convex, it is solved.

3.1 Exponential family models

We begin by defining exponential family distributions, giving several examples to illustrate a few of their properties. There are three key objects when defining a d -dimensional exponential family distribution on an underlying space \mathcal{X} : the *sufficient statistic* $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$ representing what we model, a *canonical parameter* vector $\theta \in \mathbb{R}^d$, and a *carrier* $h : \mathcal{X} \rightarrow \mathbb{R}_+$.

In the discrete case, where \mathcal{X} is a discrete set, the exponential family associated with the sufficient statistic ϕ and carrier h has probability mass function

$$p_\theta(x) = h(x) \exp(\langle \theta, \phi(x) \rangle - A(\theta)),$$

where A is the *log-partition-function*, sometimes called the *cumulant generating function*, with

$$A(\theta) := \log \sum_{x \in \mathcal{X}} h(x) \exp(\langle \theta, \phi(x) \rangle).$$

In the continuous case, p_θ is instead a density on $\mathcal{X} \subset \mathbb{R}^k$, and p_θ takes the identical form above but

$$A(\theta) = \log \int_{\mathcal{X}} h(x) \exp(\langle \theta, \phi(x) \rangle) dx.$$

We can abstract away from this distinction between discrete and continuous distributions by making the definition measure-theoretic, which we do here for completeness. (But recall the remarks in Section 1.4.)

With our notation, we have the following definition.

Definition 3.1. *The exponential family associated with the function ϕ and base measure μ is defined as the set of distributions with densities p_θ with respect to μ , where*

$$p_\theta(x) = \exp(\langle \theta, \phi(x) \rangle - A(\theta)), \quad (3.1.1)$$

and the function A is the log-partition-function (or cumulant function)

$$A(\theta) := \log \int_{\mathcal{X}} \exp(\langle \theta, \phi(x) \rangle) d\mu(x) \quad (3.1.2)$$

whenever A is finite (and is $+\infty$ otherwise). The family is regular if the domain

$$\Theta := \{\theta \mid A(\theta) < \infty\}$$

is open.

In Definition 3.1, we have included the carrier h in the base measure μ , and frequently we will give ourselves the general notation

$$p_\theta(x) = h(x) \exp(\langle \theta, \phi(x) \rangle - A(\theta)).$$

In some scenarios, it may be convenient to re-parameterize the problem in terms of some function $\eta(\theta)$ instead of θ itself; we will not worry about such issues and simply use the formulae that are most convenient.

We now give a few examples of exponential family models.

Example 3.1.1 (Bernoulli distribution): In this case, we have $X \in \{0, 1\}$ and $P(X = 1) = p$ for some $p \in [0, 1]$ in the classical version of a Bernoulli. Thus we take μ to be the counting measure on $\{0, 1\}$, and by setting $\theta = \log \frac{p}{1-p}$ to obtain a canonical representation, we have

$$\begin{aligned} P(X = x) = p(x) &= p^x(1-p)^{1-x} = \exp(x \log p - x \log(1-p)) \\ &= \exp\left(x \log \frac{p}{1-p} + \log(1-p)\right) = \exp\left(x\theta - \log(1 + e^\theta)\right). \end{aligned}$$

The Bernoulli family thus has log-partition function $A(\theta) = \log(1 + e^\theta)$. \diamond

Example 3.1.2 (Poisson distribution): The Poisson distribution (for count data) is usually parameterized by some $\lambda > 0$, and for $x \in \mathbb{N}$ has distribution $P_\lambda(X = x) = (1/x!) \lambda^x e^{-\lambda}$. Thus by taking μ to be counting (discrete) measure on $\{0, 1, \dots\}$ and setting $\theta = \log \lambda$, we find the density (probability mass function in this case)

$$p(x) = \frac{1}{x!} \lambda^x e^{-\lambda} = \exp(x \log \lambda - \lambda) \frac{1}{x!} = \exp(x\theta - e^\theta) \frac{1}{x!}.$$

Notably, taking $h(x) = (x!)^{-1}$ and log-partition $A(\theta) = e^\theta$, we have probability mass function $p_\theta(x) = h(x) \exp(\theta x - A(\theta))$. \diamond

Example 3.1.3 (Normal distribution, mean parameterization): For the d -dimensional normal distribution, we take μ to be Lebesgue measure on \mathbb{R}^d . If we fix the covariance and vary only the mean μ in the family $\mathcal{N}(\mu, \Sigma)$, then $X \sim \mathcal{N}(\mu, \Sigma)$ has density

$$p_\mu(x) = \exp \left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu) - \frac{1}{2} \log \det(2\pi\Sigma) \right).$$

Setting $h(x) = -\frac{1}{2}x^\top \Sigma^{-1}x$ and reparameterizing $\theta = \Sigma^{-1}\mu$, we obtain

$$p_\theta(x) = \underbrace{\exp \left(-\frac{1}{2}x^\top \Sigma^{-1}x - \frac{1}{2} \log \det(2\pi\Sigma) \right)}_{=:h(x)} \exp \left(x^\top \theta - \frac{1}{2}\theta^\top \Sigma \theta \right).$$

In particular, we have carrier $h(x) = \exp(-\frac{1}{2}x^\top \Sigma^{-1}x)/((2\pi)^{d/2} \det(\Sigma))$, sufficient statistic $\phi(x) = x$, and log partition $A(\theta) = \frac{1}{2}\theta^\top \Sigma^{-1}\theta$. \diamond

Example 3.1.4 (Normal distribution): Let $X \sim \mathcal{N}(\mu, \Sigma)$. We may re-parameterize this as $\Theta = \Sigma^{-1}$ and $\theta = \Sigma^{-1}\mu$, and we have density

$$p_{\theta, \Theta}(x) \propto \exp \left(\langle \theta, x \rangle - \frac{1}{2} \langle x x^\top, \Theta \rangle \right),$$

where $\langle \cdot, \cdot \rangle$ denotes the Euclidean inner product. See Exercise 3.1. \diamond

In some cases, it is analytically convenient to include a few more conditions on the exponential family.

Definition 3.2. Let $\{P_\theta\}_{\theta \in \Theta}$ be an exponential family as in Definition 3.1. The sufficient statistic ϕ is minimal if $\Theta = \text{dom } A \subset \mathbb{R}^d$ is full-dimensional and there exists no vector u such that

$$\langle u, \phi(x) \rangle \text{ is constant } \mu\text{-almost surely.}$$

Definition 3.2 is essentially equivalent to stating that $\phi(x) = (\phi_1(x), \dots, \phi_d(x))$ has linearly independent components when viewed as vectors $[\phi_i(x)]_{x \in \mathcal{X}}$. While we do not prove this, via a suitable linear transformation—a variant of Gram-Schmidt orthonormalization—one may modify any non-minimal exponential family $\{P_\theta\}$ into an equivalent minimal exponential family $\{Q_\eta\}$, meaning that the two collections satisfy the equality $\{P_\theta\} = \{Q_\eta\}$ (see Brown [41, Chapter 1]).

3.2 Why exponential families?

There are many reasons for us to study exponential families. The first major reason is their analytical tractability: as the normal distribution does, they often admit relatively straightforward computation, therefore forming a natural basis for modeling decisions. Their analytic tractability has made them the objects of substantial study for nearly the past hundred years; Brown [41] provides a deep and elegant treatment. Moreover, as we see later, they arise as the solutions to several natural optimization problems on the space of probability distributions, and they also enjoy certain robustness properties related to optimal Bayes' procedures (there is, of course, more to come on this topic).

Here, we enumerate a few of their key analytical properties, focusing on the cumulant generating (or log partition) function $A(\theta) = \log \int e^{\langle \theta, \phi(x) \rangle} d\mu(x)$. We begin with a heuristic calculation, where we assume that we exchange differentiation and integration. Assuming that this is the case, we then obtain the important expectation and covariance relationships that

$$\begin{aligned} \nabla A(\theta) &= \frac{1}{\int e^{\langle \theta, \phi(x) \rangle} d\mu(x)} \int \nabla_{\theta} e^{\langle \theta, \phi(x) \rangle} d\mu(x) \\ &= e^{-A(\theta)} \int \nabla_{\theta} e^{\langle \theta, \phi(x) \rangle} d\mu(x) = \int \phi(x) e^{\langle \theta, \phi(x) \rangle - A(\theta)} d\mu(x) = \mathbb{E}_{\theta}[\phi(X)] \end{aligned}$$

because $e^{\langle \theta, \phi(x) \rangle - A(\theta)} = p_{\theta}(x)$. A completely similar (and still heuristic, at least at this point) calculation gives

$$\nabla^2 A(\theta) = \mathbb{E}_{\theta}[\phi(X)\phi(X)^{\top}] - \mathbb{E}_{\theta}[\phi(X)]\mathbb{E}_{\theta}[\phi(X)]^{\top} = \text{Cov}_{\theta}(\phi(X)).$$

That these identities hold is no accident and is central to the appeal of exponential family models.

The first and, from our perspective, most important result about exponential family models is their convexity. While (assuming the differentiation relationships above hold) the differentiation identity that $\nabla^2 A(\theta) = \text{Cov}_{\theta}(\phi(X)) \succeq 0$ makes convexity of A immediate, one can also provide a direct argument without appealing to differentiation.

Proposition 3.2.1. *The cumulant-generating function $\theta \mapsto A(\theta)$ is convex, and it is strictly convex if and only if $\text{Cov}_{\theta}(\phi(X))$ is positive definite for all $\theta \in \text{dom } A$.*

Proof Let $\theta_{\lambda} = \lambda\theta_1 + (1-\lambda)\theta_2$, where $\theta_1, \theta_2 \in \Theta$. Then $1/\lambda \geq 1$ and $1/(1-\lambda) \geq 1$, and Hölder's inequality implies

$$\begin{aligned} \log \int \exp(\langle \theta_{\lambda}, \phi(x) \rangle) d\mu(x) &= \log \int \exp(\langle \theta_1, \phi(x) \rangle)^{\lambda} \exp(\langle \theta_2, \phi(x) \rangle)^{1-\lambda} d\mu(x) \\ &\leq \log \left(\int \exp(\langle \theta_1, \phi(x) \rangle)^{\frac{\lambda}{\lambda}} d\mu(x) \right)^{\lambda} \left(\int \exp(\langle \theta_2, \phi(x) \rangle)^{\frac{1-\lambda}{1-\lambda}} d\mu(x) \right)^{1-\lambda} \\ &= \lambda \log \int \exp(\langle \theta_1, \phi(x) \rangle) d\mu(x) + (1-\lambda) \log \int \exp(\langle \theta_2, \phi(x) \rangle) d\mu(x), \end{aligned}$$

as desired. The strict convexity will be a consequence of Proposition 3.2.2 to come, as there we formally show that $\nabla^2 A(\theta) = \text{Cov}_{\theta}(\phi(X))$. \square

We now show that $A(\theta)$ is indeed infinitely differentiable and how it generates the moments of the sufficient statistics $\phi(x)$. To describe the properties, we provide a bit of notation related to tensor products: for a vector $x \in \mathbb{R}^d$, we let

$$x^{\otimes k} := \underbrace{x \otimes x \otimes \cdots \otimes x}_{k \text{ times}}$$

denote the k th order tensor, or multilinear operator, that for $v_1, \dots, v_k \in \mathbb{R}^d$ satisfies

$$x^{\otimes k}(v_1, \dots, v_k) := \langle x, v_1 \rangle \cdots \langle x, v_k \rangle = \prod_{i=1}^k \langle x, v_i \rangle.$$

When $k = 2$, this is the familiar outer product $x^{\otimes 2} = xx^\top$. (More generally, one may think of $x^{\otimes k}$ as a $d \times d \times \cdots \times d$ box, where the (i_1, \dots, i_k) entry is $[x^{\otimes k}]_{i_1, \dots, i_k} = x_{i_1} \cdots x_{i_k}$.) With this notation, our first key result regards the differentiability of A , where we can compute (all) derivatives of $e^{A(\theta)}$ by interchanging integration and differentiation.

Proposition 3.2.2. *The cumulant-generating function $\theta \mapsto A(\theta)$ is infinitely differentiable on the interior of its domain $\Theta := \{\theta \in \mathbb{R}^d : A(\theta) < \infty\}$. The moment-generating function*

$$M(\theta) := \int \exp(\langle \theta, \phi(x) \rangle) d\mu(x)$$

is analytic on the set $\Theta_{\mathbb{C}} := \{z \in \mathbb{C}^d \mid \operatorname{Re} z \in \Theta\}$. Additionally, the derivatives of M are computed by passing through the integral, that is,

$$\begin{aligned} \nabla_{\theta}^k M(\theta) &= \nabla_{\theta}^k \int e^{\langle \theta, \phi(x) \rangle} d\mu(x) = \int \nabla_{\theta}^k e^{\langle \theta, \phi(x) \rangle} d\mu(x) \\ &= \int \phi(x)^{\otimes k} \exp(\langle \theta, \phi(x) \rangle) d\mu(x). \end{aligned}$$

The proof of the proposition is involved and requires complex analysis, so we defer it to Sec. 3.6.1.

As particular consequences of Proposition 3.2.2, we can rigorously demonstrate the expectation and covariance relationships that

$$\nabla A(\theta) = \frac{1}{\int e^{\langle \theta, \phi(x) \rangle} d\mu(x)} \int \nabla e^{\langle \theta, \phi(x) \rangle} d\mu(x) = \int \phi(x) p_{\theta}(x) d\mu(x) = \mathbb{E}_{\theta}[\phi(X)]$$

and

$$\begin{aligned} \nabla^2 A(\theta) &= \frac{1}{\int e^{\langle \theta, \phi(x) \rangle} d\mu(x)} \int \phi(x)^{\otimes 2} e^{\langle \theta, \phi(x) \rangle} d\mu(x) - \frac{(\int \phi(x) e^{\langle \theta, \phi(x) \rangle} d\mu(x))^{\otimes 2}}{(\int e^{\langle \theta, \phi(x) \rangle} d\mu(x))^2} \\ &= \mathbb{E}_{\theta}[\phi(X)\phi(X)^\top] - \mathbb{E}_{\theta}[\phi(X)]\mathbb{E}_{\theta}[\phi(X)]^\top \\ &= \operatorname{Cov}_{\theta}(\phi(X)). \end{aligned}$$

Minimal exponential families (Definition 3.2) also enjoy a few additional regularity properties. Recall that A is *strictly convex* if

$$A(\lambda\theta_0 + (1-\lambda)\theta_1) < \lambda A(\theta_0) + (1-\lambda)A(\theta_1)$$

whenever $\lambda \in (0, 1)$ and $\theta_0, \theta_1 \in \operatorname{dom} A$. We have the following proposition.

Proposition 3.2.3. *Let $\{P_{\theta}\}$ be a regular exponential family. The log partition function A is strictly convex if and only if $\{P_{\theta}\}$ is minimal.*

Proof If the family is minimal, then $\operatorname{Var}_{\theta}(u^\top \phi(X)) > 0$ for any vector u , while $\operatorname{Var}_{\theta}(u^\top \phi(X)) = u^\top \nabla^2 A(\theta) u$. This implies the strict positive definiteness $\nabla^2 A(\theta) \succ 0$, which is equivalent to strict convexity (see Corollary B.3.2 in Appendix B.3.1). Conversely, if $\nabla^2 A(\theta) \succ 0$ for all $\theta \in \Theta$, then $\operatorname{Var}_{\theta}(u^\top \phi(X)) > 0$ for all $u \neq 0$ and so $u^\top \phi(x)$ is non-constant in x . \square

3.2.1 Fitting an exponential family model

The convexity and differentiability properties make exponential family models especially attractive from a computational perspective. A major focus in statistics is the convergence of estimates of different properties of a population distribution P and whether these estimates are computable. We will develop tools to address the first of these questions, and attendant optimality guarantees, throughout this book. To set the stage for what follows, let us consider what this entails in the context of exponential family models.

Suppose we have a population P (where, for simplicity, we assume P has a density p), and for a given exponential family \mathcal{P} with densities $\{p_\theta\}$, we wish to find the model closest to P . Then it is natural (if we take on faith that the information-theoretic measures we have developed are the “right” ones) find the distribution $P_\theta \in \mathcal{P}$ closest to P in KL-divergence, that is, to solve

$$\underset{\theta}{\text{minimize}} \quad D_{\text{kl}}(P \| P_\theta) = \int p(x) \log \frac{p(x)}{p_\theta(x)} dx. \quad (3.2.1)$$

This is evidently equivalent to minimizing

$$-\int p(x) \log p_\theta(x) dx = \int p(x) [-\langle \theta, \phi(x) \rangle + A(\theta)] dx = -\langle \theta, \mathbb{E}_P[\phi(X)] \rangle + A(\theta).$$

This is always a convex optimization problem (see Appendices B and C for much more on this), as A is convex and the first term is linear, and so has no non-global optima. Here and throughout, as we mention in the introductory remarks to this chapter, we treat convex optimization as a technology: as long as the dimension of a problem is not too large and its objective can be evaluated, it is (essentially) computationally trivial.

Of course, we never have access to the population P fully; instead, we receive a sample X_1, \dots, X_n from P . In this case, a natural approach is to replace the expected (negative) log likelihood above with its empirical version and solve

$$\underset{\theta}{\text{minimize}} \quad -\sum_{i=1}^n \log p_\theta(X_i) = \sum_{i=1}^n [-\langle \theta, \phi(X_i) \rangle + A(\theta)], \quad (3.2.2)$$

which is still a convex optimization problem (as the objective is convex in θ). The maximum likelihood estimate is any vector $\hat{\theta}_n$ minimizing the negative log likelihood (3.2.2), which by setting gradients to 0 is evidently any vector satisfying

$$\nabla A(\hat{\theta}_n) = \mathbb{E}_{\hat{\theta}_n}[\phi(X)] = \frac{1}{n} \sum_{i=1}^n \phi(X_i). \quad (3.2.3)$$

In particular, we need only find a parameter $\hat{\theta}_n$ matching moments of the empirical distribution of the observed $X_i \sim P$. This $\hat{\theta}_n$ is unique whenever $\text{Cov}_\theta(\phi(X)) \succ 0$ for all θ , that is, when the covariance of ϕ is full rank in the exponential family model, because then the objective in the minimization problem (3.2.2) is strictly convex.

Let us proceed heuristically for a moment to develop a rough convergence guarantee for the estimator $\hat{\theta}_n$; the next paragraph assumes a comfort with some of classical asymptotic statistics (and the central limit theorem) and is not essential for what comes later. Then we can see how minimizers of the problem (3.2.2) converge to their population counterparts. Assume that the data

X_i are i.i.d. from an exponential family model P_{θ^*} . Then we expect that the maximum likelihood estimate $\hat{\theta}_n$ should converge to θ^* , and so

$$\frac{1}{n} \sum_{i=1}^n \phi(X_i) = \nabla A(\hat{\theta}_n) = \nabla A(\theta^*) + (\nabla^2 A(\theta^*) + o(1))(\hat{\theta}_n - \theta^*).$$

But of course, $\nabla A(\theta^*) = \mathbb{E}_{\theta^*}[\phi(X)]$, and so the central limit theorem gives that

$$\frac{1}{n} \sum_{i=1}^n (\phi(X_i) - \nabla A(\theta^*)) \dot{\sim} \mathbf{N}(0, n^{-1} \text{Cov}_{\theta^*}(\phi(X))) = \mathbf{N}(0, n^{-1} \nabla^2 A(\theta^*)),$$

where $\dot{\sim}$ means “is approximately distributed as.” Multiplying by $(\nabla^2 A(\theta^*) + o(1))^{-1} \approx \nabla^2 A(\theta^*)^{-1}$, we thus see (still working in our heuristic)

$$\begin{aligned} \hat{\theta}_n - \theta^* &= (\nabla^2 A(\theta^*) + o(1))^{-1} \frac{1}{n} \sum_{i=1}^n (\phi(X_i) - \nabla A(\theta^*)) \\ &\dot{\sim} \mathbf{N}(0, n^{-1} \cdot \nabla^2 A(\theta^*)^{-1}), \end{aligned} \tag{3.2.4}$$

where we use that $BZ \sim \mathbf{N}(0, B\Sigma B^\top)$ if $Z \sim \mathbf{N}(0, \Sigma)$. (It is possible to make each of these steps fully rigorous.) Thus the cumulant generating function A governs the error we expect in $\hat{\theta}_n - \theta^*$.

Much of the rest of this book explores properties of these types of minimization problems: at what rates do we expect $\hat{\theta}_n$ to converge to a global minimizer of problem (3.2.1)? Can we show that these rates are optimal? Is this the “right” strategy for choosing a parameter? Exponential families form a particular working example to motivate this development.

3.3 Divergence measures and information for exponential families

Their nice analytic properties mean that exponential family models also play nicely with the information theoretic tools we develop. Indeed, consider the KL-divergence between two exponential family distributions P_θ and $P_{\theta+\Delta}$, where $\Delta \in \mathbb{R}^d$. Then we have

$$\begin{aligned} D_{\text{kl}}(P_\theta \| P_{\theta+\Delta}) &= \mathbb{E}_\theta [\langle \theta, \phi(X) \rangle - A(\theta) - \langle \theta + \Delta, \phi(X) \rangle + A(\theta + \Delta)] \\ &= A(\theta + \Delta) - A(\theta) - \mathbb{E}_\theta [\langle \Delta, \phi(X) \rangle] \\ &= A(\theta + \Delta) - A(\theta) - \nabla A(\theta)^\top \Delta. \end{aligned}$$

Similarly, we have

$$\begin{aligned} D_{\text{kl}}(P_{\theta+\Delta} \| P_\theta) &= \mathbb{E}_{\theta+\Delta} [\langle \theta + \Delta, \phi(X) \rangle - A(\theta + \Delta) - \langle \theta, \phi(X) \rangle + A(\theta)] \\ &= A(\theta) - A(\theta + \Delta) + \mathbb{E}_{\theta+\Delta} [\langle \Delta, \phi(X) \rangle] \\ &= A(\theta) - A(\theta + \Delta) - \nabla A(\theta + \Delta)^\top (-\Delta). \end{aligned}$$

These identities give an immediate connection with convexity. Indeed, for a differentiable convex function h , the *Bregman divergence* associated with h is

$$D_h(u, v) = h(u) - h(v) - \langle \nabla h(v), u - v \rangle, \tag{3.3.1}$$

which is always nonnegative, and is the gap between the linear approximation to the (convex) function h and its actual value. One might more accurately call the quantity (3.3.1) the “first-order divergence,” which is more evocative, but the statistical, machine learning, and optimization literatures—in which such divergences frequently appear—have adopted this terminology, so we stick with it.

JCD Comment: Put in a picture of a Bregman divergence

We catalog these results as the following proposition.

Proposition 3.3.1. *Let $\{P_\theta\}$ be an exponential family model with cumulant generating function $A(\theta)$. Then*

$$D_{\text{kl}}(P_\theta \| P_{\theta+\Delta}) = D_A(\theta + \Delta, \theta) \quad \text{and} \quad D_{\text{kl}}(P_{\theta+\Delta} \| P_\theta) = D_A(\theta, \theta + \Delta).$$

Additionally, there exists a $t \in [0, 1]$ such that

$$D_{\text{kl}}(P_\theta \| P_{\theta+\Delta}) = \frac{1}{2} \Delta^\top \nabla^2 A(\theta + t\Delta) \Delta,$$

and similarly, there exists a $t \in [0, 1]$ such that

$$D_{\text{kl}}(P_{\theta+\Delta} \| P_\theta) = \frac{1}{2} \Delta^\top \nabla^2 A(\theta + t\Delta) \Delta.$$

Proof We have already shown the first two statements; the second two are applications of Taylor’s theorem. \square

When the perturbation Δ is small, that A is infinitely differentiable then gives that

$$D_{\text{kl}}(P_\theta \| P_{\theta+\Delta}) = \frac{1}{2} \Delta^\top \nabla^2 A(\theta) \Delta + O(\|\Delta\|^3),$$

so that the Hessian $\nabla^2 A(\theta)$ tells quite precisely how the KL divergence changes as θ varies (locally). As we saw already in Example 2.3.2 (and see the next section), when the KL-divergence between two distributions is small, it is hard to test between them, and in the sequel, we will show converses to this. The Hessian $\nabla^2 A(\theta^*)$ also governs the error in the estimate $\hat{\theta}_n - \theta^*$ in our heuristic (3.2.4). When the Hessian $\nabla^2 A(\theta)$ is quite positive semidefinite, the KL divergence $D_{\text{kl}}(P_\theta \| P_{\theta+\Delta})$ is large, and the asymptotic covariance (3.2.4) is small. For this—and other reasons we address later—for exponential family models, we call

$$\nabla^2 A(\theta) = \text{Cov}_\theta(\phi(X)) = \mathbb{E}_\theta[\nabla \log p_\theta(X) \nabla \log p_\theta(X)^\top] \quad (3.3.2)$$

the *Fisher information* of the parameter θ in the model $\{P_\theta\}$.

3.4 Generalized linear models and regression

We can specialize the general modeling strategies that exponential families provide to more directly address prediction problems, where we wish to predict a target $Y \in \mathcal{Y}$ given covariates $X \in \mathcal{X}$. Here, we almost always have that Y is either discrete or continuous with $\mathcal{Y} \subset \mathbb{R}$. In this case, we

have a sufficient statistic $\phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$, and we model $Y \mid X = x$ via the *generalized linear model* (or conditional exponential family model) if it has density or probability mass function

$$p_\theta(y \mid x) = \exp \left(\phi(x, y)^\top \theta - A(\theta \mid x) \right) h(y), \quad (3.4.1)$$

where as before h is the carrier and (in the case that $\mathcal{Y} \subset \mathbb{R}^k$)

$$A(\theta \mid x) = \log \int \exp(\phi(x, y)^\top \theta) h(y) dy$$

or, in the discrete case,

$$A(\theta \mid x) = \log \sum_y \exp(\phi(x, y)^\top \theta) h(y).$$

The log partition function $A(\cdot \mid x)$ provides the same insights for the conditional models (3.4.1) as it does for the unconditional exponential family models in the preceding sections. Indeed, as in Propositions 3.2.1 and 3.2.2, the log partition $A(\cdot \mid x)$ is always \mathcal{C}^∞ on its domain and convex. Moreover, it gives the expected moments of the sufficient statistic ϕ conditional on x , as

$$\nabla A(\theta \mid x) = \mathbb{E}_\theta[\phi(X, Y) \mid X = x],$$

and

$$\nabla^2 A(\theta \mid x) = \text{Cov}_\theta(\phi(X, Y) \mid X = x),$$

from which we can (typically) extract the mean or other statistics of Y conditional on x .

Three standard examples will be our most frequent motivators throughout this book: linear regression, binary logistic regression, and multiclass logistic regression. We give these three, as well as describing two more important examples involving modeling count data through Poisson regression and making predictions for targets y known to live in a bounded set.

Example 3.4.1 (Linear regression): In linear regression, we wish to predict $Y \in \mathbb{R}$ from a vector $X \in \mathbb{R}^d$, and assume that $Y \mid X = x$ follow the normal distribution $\mathcal{N}(\theta^\top x, \sigma^2)$. In this case, we have

$$\begin{aligned} p_\theta(y \mid x) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{1}{2\sigma^2} (y - x^\top \theta)^2 \right) \\ &= \exp \left(\frac{1}{\sigma^2} y x^\top \theta - \frac{1}{2\sigma^2} \theta^\top x x^\top \theta \right) \exp \left(-\frac{1}{2\sigma^2} y^2 + \frac{1}{2} \log(2\pi\sigma^2) \right), \end{aligned}$$

so that we have the exponential family representation (3.4.1) with $\phi(x, y) = \frac{1}{\sigma^2} x y$, $h(y) = \exp(-\frac{1}{2\sigma^2} y^2 + \frac{1}{2} \log(2\pi\sigma^2))$, and $A(\theta) = \frac{1}{2\sigma^2} \theta^\top x x^\top \theta$. As $\nabla A(\theta \mid x) = \mathbb{E}_\theta[\phi(X, Y) \mid X = x] = \frac{1}{\sigma^2} x \mathbb{E}_\theta[Y \mid X = x]$, we easily recover $\mathbb{E}_\theta[Y \mid X = x] = \theta^\top x$. \diamond

Frequently, we wish to predict binary or multiclass random variables Y . For example, consider a medical application in which we wish to assess the probability that, based on a set of covariates $x \in \mathbb{R}^d$ (say, blood pressure, height, weight, family history) and individual will have a heart attack in the next 5 years, so that $Y = 1$ indicates heart attack and $Y = -1$ indicates not. The next example shows how we might model this.

Example 3.4.2 (Binary logistic regression): If $Y \in \{-1, 1\}$, we model

$$p_\theta(y | x) = \frac{\exp(yx^\top \theta)}{1 + \exp(yx^\top \theta)},$$

where the idea in the probability above is that if $x^\top \theta$ has the same sign as y , then the large $x^\top \theta y$ becomes the higher the probability assigned the label y ; when $x^\top \theta y < 0$, the probability is small. Of course, we always have $p_\theta(y | x) + p_\theta(-y | x) = 1$, and using the identity

$$yx^\top \theta - \log(1 + \exp(yx^\top \theta)) = \frac{y+1}{2} x^\top \theta - \log(1 + \exp(x^\top \theta))$$

we obtain the generalized linear model representation $\phi(x, y) = \frac{y+1}{2}x$ and $A(\theta | x) = \log(1 + \exp(x^\top \theta))$.

As an alternative, we could represent $Y \in \{0, 1\}$ by

$$p_\theta(y | x) = \frac{\exp(yx^\top \theta)}{1 + \exp(x^\top \theta)} = \exp\left(yx^\top \theta - \log(1 + e^{x^\top \theta})\right),$$

which has the simpler sufficient statistic $\phi(x, y) = xy$. \diamond

Instead of a binary prediction problem, in many cases we have a *multiclass* prediction problem, where we seek to predict a label Y for an object x belonging to one of k different classes. For example, in image recognition, we are given an image x and wish to identify the subject Y of the image, where Y ranges over k classes, such as birds, dogs, cars, trucks, and so on. This too we can model using exponential families.

Example 3.4.3 (Multiclass logistic regression): In the case that we have a k -class prediction problem in which we wish to predict $Y \in \{1, \dots, k\}$ from $X \in \mathbb{R}^d$, we assign parameters $\theta_y \in \mathbb{R}^d$ to each of the classes $y = 1, \dots, k$. We then model

$$p_\theta(y | x) = \frac{\exp(\theta_y^\top x)}{\sum_{j=1}^k \exp(\theta_j^\top x)} = \exp\left(\theta_y^\top x - \log\left(\sum_{j=1}^k e^{\theta_j^\top x}\right)\right).$$

Here, the idea is that if $\theta_y^\top x > \theta_j^\top x$ for all $j \neq y$, then the model assigns higher probability to class y than any other class; the larger the gap between $\theta_y^\top x$ and $\theta_j^\top x$, the larger the difference in assigned probabilities. \diamond

Other approaches with these ideas allow us to model other situations. Poisson regression models are frequent choices for modeling count data. For example, consider an insurance company that wishes to issue premiums for shipping cargo in different seasons and on different routes, and so wishes to predict the number of times a given cargo ship will be damaged by waves over a period of service; we might represent this with a feature vector x encoding information about the ship to be insured, typical weather on the route it will take, and the length of time it will be in service. To model such counts $Y \in \{0, 1, 2, \dots\}$, we turn to Poisson regression.

Example 3.4.4 (Poisson regression): When $Y \in \mathbb{N}$ is a count, the Poisson distribution with rate $\lambda > 0$ gives $P(Y = y) = \frac{e^{-\lambda} \lambda^y}{y!}$. Poisson regression models λ via $e^{\theta^\top x}$, giving model

$$p_\theta(y | x) = \frac{1}{y!} \exp\left(yx^\top \theta - e^{\theta^\top x}\right),$$

so that we have carrier $h(y) = 1/y!$ and the simple sufficient statistic $yx^\top \theta$. The log partition function is $A(\theta | x) = e^{\theta^\top x}$. \diamond

Lastly, we consider a less standard example, but which highlights the flexibility of these models. Here, we assume a linear regression problem but in which we wish to predict values Y in a bounded range.

Example 3.4.5 (Bounded range regression): Suppose that we know $Y \in [-b, b]$, but we wish to model it via an exponential family model with density

$$p_\theta(y | x) = \exp(yx^\top \theta - A(\theta | x)) \mathbf{1}\{y \in [-b, b]\},$$

which is non-zero only for $-b \leq y \leq b$. Letting $s = x^\top \theta$ for shorthand, we have

$$\int_{-b}^b e^{ys} dy = \frac{1}{s} [e^{bs} - e^{-bs}],$$

where the limit as $s \rightarrow 0$ is $2b$; the (conditional) log partition function is thus

$$A(\theta | x) = \begin{cases} \log \frac{e^{b\theta^\top x} - e^{-b\theta^\top x}}{\theta^\top x} & \text{if } \theta^\top x \neq 0 \\ \log(2b) & \text{otherwise.} \end{cases}$$

While its functional form makes this highly non-obvious, our general results guarantee that $A(\theta | x)$ is indeed \mathcal{C}^∞ and convex in θ . We have $\nabla A(\theta | x) = x\mathbb{E}_\theta[Y | X = x]$ because $\phi(x, y) = xy$, and we can therefore immediately recover $\mathbb{E}_\theta[Y | X = x]$. Indeed, set $s = \theta^\top x$, and without loss of generality assume $s \neq 0$. Then

$$\mathbb{E}[Y | x^\top \theta = s] = \frac{\partial}{\partial s} \log \frac{e^{bs} - e^{-bs}}{s} = \frac{b(e^{bs} + e^{-bs})}{e^{bs} - e^{-bs}} - \frac{1}{s},$$

which increases from $-b$ to b as $s = x^\top \theta$ increases from $-\infty$ to $+\infty$. \diamond

3.4.1 Fitting a generalized linear model from a sample

We briefly revisit the approach in Section 3.2.1 for fitting exponential family models in the context of generalized linear models. In this case, the analogue of the maximum likelihood problem (3.2.2) is to solve

$$\underset{\theta}{\text{minimize}} \quad -\sum_{i=1}^n \log p_\theta(Y_i | X_i) = \sum_{i=1}^n [-\phi(X_i, Y_i)^\top \theta + A(\theta | X_i)].$$

This is a convex optimization problem with \mathcal{C}^∞ objective, so we can treat solving it as an (essentially) trivial problem unless the sample size n or dimension d of θ are astronomically large.

As in the moment matching equality (3.2.3), a necessary and sufficient condition for $\hat{\theta}_n$ to minimize the above objective is that it achieves 0 gradient, that is,

$$\frac{1}{n} \sum_{i=1}^n \nabla A(\hat{\theta}_n | X_i) = \frac{1}{n} \sum_{i=1}^n \phi(X_i, Y_i).$$

Once again, to find $\hat{\theta}_n$ amounts to matching moments, as $\nabla A(\theta | X_i) = \mathbb{E}[\phi(X, Y) | X = X_i]$, and we still enjoy the convexity properties of the standard exponential family models.

In general, we of course do not expect any exponential family or generalized linear model (GLM) to have perfect fidelity to the world: all models are inaccurate (but many are useful!). Nonetheless,

we can still *fit* any of the GLM models in Examples 3.4.1–3.4.5 to data of the appropriate type. In particular, for the logarithmic loss $\ell(\theta; x, y) = -\log p_\theta(y | x)$, we can define the empirical loss

$$L_n(\theta) := \frac{1}{n} \sum_{i=1}^n \ell(\theta; X_i, Y_i).$$

Then, as $n \rightarrow \infty$, we expect that $L_n(\theta) \rightarrow \mathbb{E}[\ell(\theta; X, Y)]$, so that the minimizing θ should give the best predictions possible according to the loss ℓ . We shall therefore often be interested in such convergence guarantees and the deviations of sample quantities (like L_n) from their population counterparts.

3.4.2 The information in a generalized linear model

As we did in Section 3.3, we can compute the “information” about a parameter θ in a generalized linear model as well. In this case, P_θ specifies only the conditional distribution of Y given $X = x$, so when we compute the information, we assume X follows some marginal distribution Q . In this case, $(X, Y) \sim P_\theta \circ Q$, where we have abused composition notation to mean that $\mathbb{P}(X, Y \in A) = \int_{(x,y) \in A} p_\theta(y | x) q(x) dy dx$. In this case, we have

$$\begin{aligned} D_{\text{kl}}(P_\theta \circ Q \| P_{\theta+\Delta} \circ Q) &= \mathbb{E}_\theta \left[\log \frac{p_\theta(Y | X)}{p_{\theta+\Delta}(Y | X)} \right] = \mathbb{E}_\theta \left[\Delta^\top \phi(X, Y) - A(\theta | X) + A(\theta + \Delta | X) \right] \\ &= \mathbb{E}_\theta [D_{A(\cdot|X)}(\theta + \Delta, \theta)], \end{aligned}$$

and similarly

$$D_{\text{kl}}(P_{\theta+\Delta} \circ Q \| P_\theta \circ Q) = \mathbb{E}_{\theta+\Delta} [D_{A(\cdot|X)}(\theta, \theta + \Delta)],$$

where we recall the Bregman divergence (3.3.1) and have used that $\mathbb{E}_\theta[\phi(X, Y) | X] = \nabla A(\theta | X)$. Performing a Taylor expansion, we have

$$A(\theta + \Delta | x) = A(\theta | x) + \langle \nabla A(\theta | x), \Delta \rangle + \frac{1}{2} \Delta^\top \nabla^2 A(\theta | x) \Delta + O(\mathbb{E}_{\theta+t\Delta}[\|\phi(X, Y)\|^3 | x] \cdot \|\Delta\|^3),$$

where we have computed third derivatives of $A(\theta | x)$, and $t \in [0, 1]$. Evaluating the Taylor expansion in the integral form, once there exists some $\delta > 0$ such that

$$\int_0^1 \mathbb{E}_Q [\mathbb{E}_{\theta+tv}[\|\phi(X, Y)\|^3 | X]] dt < \infty$$

for any $\|v\|_2 \leq \delta$, for either of the expansions above we have the following corollary:

Corollary 3.4.6. *Assume the marginal distribution Q on X satisfies the above integrability condition. Then as $\Delta \rightarrow 0$,*

$$D_{\text{kl}}(P_\theta \circ Q \| P_{\theta+\Delta} \circ Q) = \frac{1}{2} \Delta^\top \mathbb{E}_Q [\nabla^2 A(\theta | X)] \Delta + O(\|\Delta\|^3)$$

and

$$D_{\text{kl}}(P_{\theta+\Delta} \circ Q \| P_\theta \circ Q) = \frac{1}{2} \Delta^\top \mathbb{E}_Q [\nabla^2 A(\theta | X)] \Delta + O(\|\Delta\|^3).$$

In analogy with Proposition 3.3.1, we see again that the expected Hessian $\mathbb{E}_Q[\nabla^2 A(\theta | X)]$ tells quite precisely how the KL divergence changes as θ varies locally, but now, the distribution Q on X also enters the picture. So when A and the distribution Q are such that $\mathbb{E}_Q[\nabla^2 A(\theta | X)]$ is large in the semidefinite order, then it is easy to distinguish data coming from $P_\theta \circ Q$ from that drawn from $P_{\theta'} \circ Q$, and otherwise, it is not. We therefore call

$$\mathbb{E}[\nabla^2 A(\theta | X)] = \mathbb{E}[\text{Cov}_\theta(\phi(X, Y) | X)] = \mathbb{E}[\nabla \log p_\theta(Y | X) \nabla \log p_\theta(Y | X)] \quad (3.4.2)$$

the *Fisher information* of the parameter θ in the model P_θ .

Example 3.4.7 (The information in logistic regression): For the binary logistic regression model (Example 3.4.2) with $Y \in \{0, 1\}$, we have

$$\nabla \log p_\theta(y | x) = yx - \frac{e^{x^\top \theta}}{1 + e^{x^\top \theta}} x = (y - p_\theta(1 | x))x$$

and $\nabla^2 A(\theta | X) = p_\theta(1 | x)(1 - p_\theta(1 | x))xx^\top$. Thus for $X \sim Q$ we the Fisher information is

$$\mathbb{E}_Q[\text{Var}_\theta(Y | X)XX^\top] \preceq \frac{1}{4}\mathbb{E}_Q[XX^\top].$$

When $\theta = 0$, we have identically $\text{Var}_\theta(Y | X) = \frac{1}{4}$, which is the “maximal” information. Additionally, we see that when $X \sim Q$ has larger covariance, we expect XX^\top to be larger in the semidefinite order, meaning the observations (X, Y) contain more information about θ (of course, this is mitigated by the fact that $p_\theta(y | x)$ becomes more extreme as $\|x\|$ grows). \diamond

Example 3.4.8 (The KL-divergence in logistic regression): The binary logistic regression model with $Y \in \{0, 1\}$ also admits simple bounds on its KL-divergence. For these, we first make the simple observation that for the log-sum-exp function $f(t) = \log(1 + e^t)$, we have $f'(t) = \frac{e^t}{1+e^t}$ and $f''(t) = f'(t)(1 - f'(t))$. Taylor’s theorem states that

$$f(t + \Delta) = f(t) + f'(t)\Delta + \int_0^\Delta f''(t + u)(\Delta - u)du,$$

and as $0 \leq f'' \leq \frac{1}{4}$, we have $|\int_0^\Delta f''(t + u)(\Delta - u)du| \leq \frac{1}{4} \int_0^\Delta |\Delta - u|du = \frac{\Delta^2}{8}$, so

$$|f(t + \Delta) - f(t) - f'(t)\Delta| \leq \frac{\Delta^2}{8}$$

for all $\Delta \in \mathbb{R}$. Computing the KL-divergence directly, we thus have for any parameters θ_0, θ_1 that

$$D_{\text{kl}}(P_{\theta_0}(\cdot | x) \| P_{\theta_1}(\cdot | x)) = f(\theta_0^\top x) - f(\theta_1^\top x) - f'(\theta_1^\top x)x^\top(\theta_0 - \theta_1) \leq \frac{1}{8} \left(x^\top(\theta_0 - \theta_1)\right)^2,$$

so the divergence is at most quadratic. \diamond

3.5 Lower bounds on testing a parameter's value

We give a bit of a preview here of the tools we will develop to prove fundamental limits in Part II of the book, an *hors d'oeuvres* that points to the techniques we develop. In Section 2.3.1, we presented Le Cam's method and used it in Example 2.3.2 to give a lower bound on the probability of error in a hypothesis test comparing two normal means. This approach extends beyond this simple case, and here we give another example applying it to exponential family models.

We give a stylized version of the problem. Let $\{P_\theta\}$ be an exponential family model with parameter $\theta \in \mathbb{R}^d$. Suppose for some vector $v \in \mathbb{R}^d$, we wish to test whether $v^\top \theta > 0$ or $v^\top \theta < 0$ in the model. For example, in the regression settings in Section 3.4, we may be interested in the effect of a treatment on health outcomes. Then the covariates x contain information about an individual with first index x_1 corresponding to whether the individual is treated or not, while Y measures the outcome of treatment; setting $v = e_1$, we then wish to test whether there is a positive treatment effect $\theta_1 = e_1^\top \theta > 0$ or negative.

Abstracting away the specifics of the scenario, we ask the following question: given an exponential family $\{P_\theta\}$ and a threshold t of interest, at what separation $\delta > 0$ does it become essentially impossible to test

$$v^\top \theta \leq t \quad \text{versus} \quad v^\top \theta \geq t + \delta?$$

We give one approach to this using two-point hypothesis testing lower bounds. In this case, we consider testing sequences of two alternatives

$$H_0 : \theta = \theta_0 \quad \text{versus} \quad H_{1,n} : \theta = \theta_n$$

as n grows, where we observe a sample X_1^n drawn i.i.d. either according to P_{θ_0} (i.e., H_0) or P_{θ_n} (i.e., $H_{1,n}$). By choosing θ_n in a way that makes the separation $v^\top (\theta_n - \theta_0)$ large but testing H_0 against $H_{1,n}$ challenging, we can then (roughly) identify the separation δ at which testing becomes impossible.

Proposition 3.5.1. *Let $\theta_0 \in \mathbb{R}^d$. Then there exists a sequence of parameters θ_n with $\|\theta_n - \theta_0\| = O(1/\sqrt{n})$, separation*

$$v^\top (\theta_n - \theta_0) = \frac{1}{\sqrt{n}} \sqrt{v^\top \nabla^2 A(\theta_0)^{-1} v},$$

and for which

$$\inf_{\Psi} \{P_{\theta_0}(\Psi(X_1^n) \neq 0) + P_{\theta_n}(\Psi(X_1^n) \neq 1)\} \geq \frac{1}{2} + O(n^{-1/2}).$$

Proof Let $\Delta \in \mathbb{R}^d$ be a potential perturbation to $\theta_1 = \theta_0 + \Delta$, which gives separation $\delta = v^\top \theta_1 - v^\top \theta_0 = v^\top \Delta$. Let $P_0 = P_{\theta_0}$ and $P_1 = P_{\theta_1}$. Then the smallest summed probability of error in testing between P_0 and P_1 based on n observations X_1^n is

$$\inf_{\Psi} \{P_0(\Psi(X_1, \dots, X_n) \neq 0) + P_1(\Psi(X_1, \dots, X_n) \neq 1)\} = 1 - \|P_0^n - P_1^n\|_{\text{TV}}$$

by Proposition 2.3.1. Following the approach of Example 2.3.2, we apply Pinsker's inequality (2.2.10) and use that the KL-divergence tensorizes to find

$$2 \|P_0^n - P_1^n\|_{\text{TV}}^2 \leq n D_{\text{kl}}(P_0 \| P_1) = n D_{\text{kl}}(P_{\theta_0} \| P_{\theta_0 + \Delta}) = n D_A(\theta_0 + \Delta, \theta_0),$$

where the final equality follows from the equivalence between KL and Bregman divergences for exponential families (Proposition 3.3.1).

To guarantee that the summed probability of error is at least $\frac{1}{2}$, that is, $\|P_0^n - P_1^n\|_{\text{TV}} \leq \frac{1}{2}$, it suffices to choose Δ satisfying $nD_A(\theta_0 + \Delta, \theta_0) \leq \frac{1}{2}$. So to maximize the separation $v^\top \Delta$ while guaranteeing a constant probability of error, we (approximately) solve

$$\begin{aligned} & \text{maximize} && v^\top \Delta \\ & \text{subject to} && D_A(\theta_0 + \Delta, \theta_0) \leq \frac{1}{2n}. \end{aligned}$$

Now, consider that $D_A(\theta_0 + \Delta, \theta_0) = \frac{1}{2} \Delta^\top \nabla^2 A(\theta_0) \Delta + O(\|\Delta\|^3)$. Ignoring the higher order term, we consider maximizing $v^\top \Delta$ subject to $\Delta^\top \nabla^2 A(\theta_0) \Delta \leq \frac{1}{n}$. A Lagrangian calculation shows that this has solution

$$\Delta = \frac{1}{\sqrt{n}} \frac{1}{\sqrt{v^\top \nabla^2 A(\theta_0)^{-1} v}} \nabla^2 A(\theta_0)^{-1} v.$$

With this choice, we have separation $\delta = v^\top \Delta = \sqrt{v^\top \nabla^2 A(\theta_0)^{-1} v / n}$, and $D_A(\theta_0 + \Delta, \theta_0) = \frac{1}{2n} + O(1/n^{3/2})$. The summed probability of error is at least

$$1 - \|P_0^n - P_1^n\|_{\text{TV}} \geq 1 - \sqrt{\frac{n}{4n} + O(n^{-1/2})} = 1 - \sqrt{\frac{1}{4} + O(n^{-1/2})} = \frac{1}{2} + O(n^{-1/2})$$

as desired. \square

Let us briefly sketch out why Proposition 3.5.1 is the “right” answer using the heuristics in Section 3.2.1. For an unknown parameter θ in the exponential family model P_θ , we observe X_1, \dots, X_n , and wish to test whether $v^\top \theta \geq t$ for a given threshold t . Call our null $H_0 : v^\top \theta \leq t$, and assume we wish to test at an asymptotic level $\alpha > 0$, meaning the probability the test falsely rejects H_0 is (as $n \rightarrow \infty$) is at most α . Assuming the heuristic (3.2.4), we have the approximate distributional equality

$$v^\top \hat{\theta}_n \sim \mathcal{N}\left(v^\top \theta, \frac{1}{n} v^\top \nabla^2 A(\hat{\theta}_n)^{-1} v\right).$$

Note that we have $\hat{\theta}_n$ on the right side of the distribution; it is possible to make this rigorous, but here we target only intuition building. A natural asymptotically level α test is then

$$T_n := \begin{cases} \text{Reject} & \text{if } v^\top \hat{\theta}_n \geq t + z_{1-\alpha} \sqrt{v^\top \nabla^2 A(\hat{\theta}_n)^{-1} v / n} \\ \text{Accept} & \text{otherwise,} \end{cases}$$

where $z_{1-\alpha}$ is the $1 - \alpha$ quantile of a standard normal, $\mathbb{P}(Z \geq z_{1-\alpha}) = \alpha$ for $Z \sim \mathcal{N}(0, 1)$. Let θ_0 be such that $v^\top \theta_0 = t$, so H_0 holds. Then

$$P_{\theta_0}(T_n \text{ rejects}) = P_{\theta_0}\left(\sqrt{n} \cdot v^\top (\hat{\theta}_n - \theta_0) \geq z_{1-\alpha} \sqrt{v^\top \nabla^2 A(\hat{\theta}_n)^{-1} v}\right) \rightarrow \alpha.$$

At least heuristically, then, this separation $\delta = \sqrt{v^\top \nabla^2 A(\theta_0)^{-1} v} / \sqrt{n}$ is the fundamental separation in parameter values at which testing becomes possible (or below which it is impossible).

As a brief and suggestive aside, the precise growth of the KL-divergence $D_{\text{kl}}(P_{\theta_0 + \Delta} \| P_{\theta_0}) = \frac{1}{2} \Delta^\top \nabla^2 A(\theta_0) \Delta + O(\|\Delta\|^3)$ near θ_0 plays the fundamental role in both the lower bound and upper bound on testing. When the Hessian $\nabla^2 A(\theta_0)$ is “large,” meaning it is very positive definite, distributions with small parameter distances are still well-separated in KL-divergence, making testing easy, while when $\nabla^2 A(\theta_0)$ is small (nearly indefinite), the KL-divergence can be small even for large parameter separations Δ and testing is hard. As a consequence, at least for exponential family models, the Fisher information (3.3.2), which we defined as $\nabla^2 A(\theta) = \text{Cov}_\theta(\phi(X))$, plays a central role in testing and, as we see later, estimation.

3.6 Deferred proofs

We collect proofs that rely on background we do not assume for this book here.

3.6.1 Proof of Proposition 3.2.2

We follow Brown [41]. We demonstrate only the first-order differentiability using Lebesgue's dominated convergence theorem, as higher orders and the interchange of integration and differentiation are essentially identical. Demonstrating first-order complex differentiability is of course enough to show that A is analytic.¹ As the proof of Proposition 3.2.1 does not rely on analyticity of A , we may use its results. Thus, let $\Theta = \text{dom } A(\cdot)$ in \mathbb{R}^d , which is convex. We assume Θ has non-empty interior (if the interior is empty, then the convexity of Θ means that it must lie in a lower dimensional subspace; we simply take the interior relative to that subspace and may proceed). We claim the following lemma, which is the key to applying dominated convergence; we state it first for \mathbb{R}^d .

Lemma 3.6.1. *Consider any collection $\{\theta_1, \dots, \theta_m\} \subset \Theta$, and let $\Theta_0 = \text{Conv}\{\theta_i\}_{i=1}^m$ and $C \subset \text{int } \Theta_0$. Then for any $k \in \mathbb{N}$, there exists a constant $K = K(C, k, \{\theta_i\})$ such that for all $\theta_0 \in C$,*

$$\|x\|^k \exp(\langle \theta_0, x \rangle) \leq K \max_{j \leq m} \exp(\langle \theta_j, x \rangle).$$

Proof Let $\mathbb{B} = \{u \in \mathbb{R}^d \mid \|u\| \leq 1\}$ be the unit ball in \mathbb{R}^d . For any $\epsilon > 0$, there exists a $K = K(\epsilon)$ such that $\|x\|^k \leq K e^{\epsilon \|x\|}$ for all $x \in \mathbb{R}^d$. As $C \subset \text{int Conv}(\Theta_0)$, there exists an $\epsilon > 0$ such that for all $\theta_0 \in C$, $\theta_0 + 2\epsilon \mathbb{B} \subset \Theta_0$, and by construction, for any $u \in \mathbb{B}$ we can write $\theta_0 + 2\epsilon u = \sum_{j=1}^m \lambda_j \theta_j$ for some $\lambda \in \mathbb{R}_+^m$ with $\mathbf{1}^\top \lambda = 1$. We therefore have

$$\begin{aligned} \|x\|^k \exp(\langle \theta_0, x \rangle) &\leq \|x\|^k \sup_{u \in \mathbb{B}} \exp(\langle \theta_0 + \epsilon u, x \rangle) \\ &= \|x\|^k \exp(\epsilon \|x\|) \exp(\langle \theta_0, x \rangle) \leq K \exp(2\epsilon \|x\|) \exp(\langle \theta_0, x \rangle) \\ &= K \sup_{u \in \mathbb{B}} \exp(\langle \theta_0 + 2\epsilon u, x \rangle). \end{aligned}$$

But using the convexity of $t \mapsto \exp(t)$ and that $\theta_0 + 2\epsilon u \in \Theta_0$, the last quantity has upper bound

$$\sup_{u \in \mathbb{B}} \exp(\langle \theta_0 + 2\epsilon u, x \rangle) \leq \max_{j \leq m} \exp(\langle \theta_j, x \rangle).$$

This gives the desired claim. □

A similar result is possible with differences of exponentials:

Lemma 3.6.2. *Under the conditions of Lemma 3.6.1, there exists a K such that for any $\theta, \theta_0 \in C$*

$$\frac{e^{\langle \theta, x \rangle} - e^{\langle \theta_0, x \rangle}}{\|\theta - \theta_0\|} \leq K \max_{j \leq m} e^{\langle \theta_j, x \rangle}.$$

Proof We write

$$\frac{\exp(\langle \theta, x \rangle) - \exp(\langle \theta_0, x \rangle)}{\|\theta - \theta_0\|} = \frac{\exp(\langle \theta - \theta_0, x \rangle) - 1}{\|\theta - \theta_0\|} \exp(\langle \theta_0, x \rangle)$$

¹For complex functions, Osgood's lemma shows that if A is continuous and holomorphic in each variable individually, it is holomorphic. For a treatment of such ideas in an engineering context, see, e.g. [101, Ch. 1].

so that the lemma is equivalent to showing that

$$\frac{|e^{\langle \theta - \theta_0, x \rangle} - 1|}{\|\theta - \theta_0\|} \leq K \max_{j \leq m} \exp(\langle \theta_j - \theta_0, x \rangle).$$

From this, we can assume without loss of generality that $\theta_0 = \mathbf{0}$ (by shifting). Now note that by convexity $e^{-a} \geq 1 - a$ for all $a \in \mathbb{R}$, so $1 - e^a \leq |a|$ when $a \leq 0$. Conversely, if $a > 0$, then $ae^a \geq e^a - 1$ (note that $\frac{d}{da}(ae^a) = ae^a + e^a \geq e^a$), so dividing by $\|x\|$, we see that

$$\frac{|e^{\langle \theta, x \rangle} - 1|}{\|\theta\| \|x\|} \leq \frac{|e^{\langle \theta, x \rangle} - 1|}{|\langle \theta, x \rangle|} \leq \frac{\max\{\langle \theta, x \rangle e^{\langle \theta, x \rangle}, |\langle \theta, x \rangle|\}}{|\langle \theta, x \rangle|} \leq e^{\langle \theta, x \rangle} + 1.$$

As $\theta \in C$, Lemma 3.6.1 then implies that

$$\frac{|e^{\langle \theta, x \rangle} - 1|}{\|\theta\|} \leq \|x\| (e^{\langle \theta, x \rangle} + 1) \leq K \max_j e^{\langle \theta_j, x \rangle},$$

as desired. \square

With the lemmas in hand, we can demonstrate a dominating function for the derivatives. Indeed, fix $\theta_0 \in \text{int } \Theta$ and for $\theta \in \Theta$, define

$$g(\theta, x) = \frac{\exp(\langle \theta, x \rangle) - \exp(\langle \theta_0, x \rangle) - \exp(\langle \theta_0, x \rangle) \langle x, \theta - \theta_0 \rangle}{\|\theta - \theta_0\|} = \frac{e^{\langle \theta, x \rangle} - e^{\langle \theta_0, x \rangle} - \langle \nabla e^{\langle \theta_0, x \rangle}, \theta - \theta_0 \rangle}{\|\theta - \theta_0\|}.$$

Then $\lim_{\theta \rightarrow \theta_0} g(\theta, x) = 0$ by the differentiability of $t \mapsto e^t$. Lemmas 3.6.1 and 3.6.2 show that if we take any collection $\{\theta_j\}_{j=1}^m \subset \Theta$ for which $\theta \in \text{int Conv}\{\theta_j\}$, then for $C \subset \text{int Conv}\{\theta_j\}$, there exists a constant K such that

$$|g(\theta, x)| \leq \frac{|\exp(\langle \theta, x \rangle) - \exp(\langle \theta_0, x \rangle)|}{\|\theta - \theta_0\|} + \|x\| \exp(\langle \theta_0, x \rangle) \leq K \max_j \exp(\langle \theta_j, x \rangle)$$

for all $\theta \in C$. As $\int \max_j e^{\langle \theta_j, x \rangle} d\mu(x) \leq \sum_{j=1}^m \int e^{\langle \theta_j, x \rangle} d\mu(x) < \infty$, the dominated convergence theorem thus implies that

$$\lim_{\theta \rightarrow \theta_0} \int g(\theta, x) d\mu(x) = 0,$$

and so $M(\theta) = \exp(A(\theta))$ is differentiable in θ , as

$$M(\theta) = M(\theta_0) + \left\langle \int x e^{\langle \theta_0, x \rangle} d\mu(x), \theta - \theta_0 \right\rangle + o(\|\theta - \theta_0\|).$$

It is evident that we have the derivative

$$\nabla M(\theta) = \int \nabla \exp(\langle \theta, x \rangle) d\mu(x).$$

Analyticity Over the subset $\Theta_{\mathbb{C}} := \{\theta + iz \mid \theta \in \Theta, z \in \mathbb{R}^d\}$ (where $i = \sqrt{-1}$ is the imaginary unit), we can extend the preceding results to demonstrate that A is analytic on $\Theta_{\mathbb{C}}$. Indeed, we first simply note that for $a, b \in \mathbb{R}$, $\exp(a + ib) = \exp(a) \exp(ib)$ and $|\exp(a + ib)| = \exp(a)$, i.e. $|e^z| = e^{\text{Re } z}$ for $z \in \mathbb{C}$, and so Lemmas 3.6.1 and 3.6.2 follow *mutatis-mutandis* as in the real case. These are enough for the application of the dominated convergence theorem above, and we use that $\exp(\cdot)$ is analytic to conclude that $\theta \mapsto M(\theta)$ is analytic on $\Theta_{\mathbb{C}}$.

3.7 Bibliography

3.8 Exercises

Exercise 3.1: In Example 3.1.4, give the sufficient statistic ϕ and an explicit formula for the log partition function $A(\theta, \Theta)$ so that we can write $p_{\theta, \Theta}(x) = \exp(\langle \theta, \phi_1(x) \rangle + \langle \Theta, \phi_2(x) \rangle - A(\theta, \Theta))$.

Exercise 3.2: Consider the binary logistic regression model in Example 3.4.2, and let $\ell(\theta; x, y) = -\log p_{\theta}(y | x)$ be the associated log loss.

(i) Give the Hessian $\nabla_{\theta}^2 \ell(\theta; x, y)$.

(ii) Let $(x_i, y_i)_{i=1}^n \subset \mathbb{R}^d \times \{\pm 1\}$ be a sample. Give a sufficient condition for the minimizer of the empirical log loss

$$L_n(\theta) := \frac{1}{n} \sum_{i=1}^n \ell(\theta; x_i, y_i)$$

to be unique that depends only on the vectors $\{x_i\}$. *Hint.* A convex function h is strictly convex if and only if its Hessian $\nabla^2 h$ is positive definite.

Exercise 3.3: Give the Fisher information (3.4.2) for each of the following generalized linear models:

(a) Linear regression (Example 3.4.1).

(b) Poisson regression (Example 3.4.4).

Part I

Concentration, information, stability, and generalization

Chapter 4

Concentration Inequalities

In many scenarios, it is useful to understand how a random variable X behaves by giving bounds on the probability that it deviates far from its mean or median. This can allow us to give prove that estimation and learning procedures will have certain performance, that different decoding and encoding schemes work with high probability, among other results. In this chapter, we give several tools for proving bounds on the probability that random variables are far from their typical values. We conclude the section with a discussion of basic uniform laws of large numbers and applications to empirical risk minimization and statistical learning, though we focus on the relatively simple cases we can treat with our tools.

4.1 Basic tail inequalities

In this first section, we have a simple to state goal: given a random variable X , how does X concentrate around its mean? That is, assuming w.l.o.g. that $\mathbb{E}[X] = 0$, how well can we bound

$$\mathbb{P}(X \geq t)?$$

We begin with the three most classical three inequalities for this purpose: the Markov, Chebyshev, and Chernoff bounds, which are all instances of the same technique.

The basic inequality off of which all else builds is Markov's inequality.

Proposition 4.1.1 (Markov's inequality). *Let X be a nonnegative random variable, meaning that $X \geq 0$ with probability 1. Then*

$$\mathbb{P}(X \geq t) \leq \frac{\mathbb{E}[X]}{t}.$$

Proof For any random variable, $\mathbb{P}(X \geq t) = \mathbb{E}[\mathbf{1}\{X \geq t\}] \leq \mathbb{E}[(X/t)\mathbf{1}\{X \geq t\}] \leq \mathbb{E}[X]/t$, as $X/t \geq 1$ whenever $X \geq t$. \square

When we know more about a random variable than that its expectation is finite, we can give somewhat more powerful bounds on the probability that the random variable deviates from its typical values. The first step in this direction, Chebyshev's inequality, requires two moments, and when we have exponential moments, we can give even stronger results. As we shall see, each of these results is but an application of Proposition 4.1.1.

Proposition 4.1.2 (Chebyshev's inequality). *Let X be a random variable with $\text{Var}(X) < \infty$. Then*

$$\mathbb{P}(X - \mathbb{E}[X] \geq t) \leq \frac{\text{Var}(X)}{t^2} \quad \text{and} \quad \mathbb{P}(X - \mathbb{E}[X] \leq -t) \leq \frac{\text{Var}(X)}{t^2}$$

for all $t \geq 0$.

Proof We prove only the upper tail result, as the lower tail is identical. We first note that $X - \mathbb{E}[X] \geq t$ implies that $(X - \mathbb{E}[X])^2 \geq t^2$. But of course, the random variable $Z = (X - \mathbb{E}[X])^2$ is nonnegative, so Markov's inequality gives $\mathbb{P}(X - \mathbb{E}[X] \geq t) \leq \mathbb{P}(Z \geq t^2) \leq \mathbb{E}[Z]/t^2$, and $\mathbb{E}[Z] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \text{Var}(X)$. \square

If a random variable has a moment generating function—exponential moments—we can give bounds that enjoy very nice properties when combined with sums of random variables. First, we recall that

$$\varphi_X(\lambda) := \mathbb{E}[e^{\lambda X}]$$

is the moment generating function of the random variable X . Then we have the Chernoff bound.

Proposition 4.1.3. *For any random variable X , we have*

$$\mathbb{P}(X \geq t) \leq \frac{\mathbb{E}[e^{\lambda X}]}{e^{\lambda t}} = \varphi_X(\lambda) e^{-\lambda t}$$

for all $\lambda \geq 0$.

Proof This is another application of Markov's inequality: for $\lambda > 0$, we have $e^{\lambda X} \geq e^{\lambda t}$ if and only if $X \geq t$, so that $\mathbb{P}(X \geq t) = \mathbb{P}(e^{\lambda X} \geq e^{\lambda t}) \leq \mathbb{E}[e^{\lambda X}]/e^{\lambda t}$. \square

In particular, taking the infimum over all $\lambda \geq 0$ in Proposition 4.1.3 gives the more standard Chernoff (large deviation) bound

$$\mathbb{P}(X \geq t) \leq \exp \left(\inf_{\lambda \geq 0} \log \varphi_X(\lambda) - \lambda t \right).$$

Example 4.1.4 (Gaussian random variables): When X is a mean-zero Gaussian variable with variance σ^2 , we have

$$\varphi_X(\lambda) = \mathbb{E}[\exp(\lambda X)] = \exp \left(\frac{\lambda^2 \sigma^2}{2} \right). \quad (4.1.1)$$

To see this, we compute the integral; we have

$$\begin{aligned} \mathbb{E}[\exp(\lambda X)] &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(\lambda x - \frac{1}{2\sigma^2} x^2 \right) dx \\ &= e^{\frac{\lambda^2 \sigma^2}{2}} \underbrace{\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{1}{2\sigma^2} (x - \lambda \sigma^2)^2 \right) dx}_{=1} \end{aligned}$$

because this is simply the integral of the Gaussian density.

As a consequence of the equality (4.1.1) and the Chernoff bound technique (Proposition 4.1.3), we see that for X Gaussian with variance σ^2 , we have

$$\mathbb{P}(X \geq \mathbb{E}[X] + t) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right) \quad \text{and} \quad \mathbb{P}(X \leq \mathbb{E}[X] - t) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right)$$

for all $t \geq 0$. Indeed, we have $\log \varphi_{X-\mathbb{E}[X]}(\lambda) = \frac{\lambda^2 \sigma^2}{2}$, and $\inf_{\lambda} \{\frac{\lambda^2 \sigma^2}{2} - \lambda t\} = -\frac{t^2}{2\sigma^2}$, which is attained by $\lambda = \frac{t}{\sigma^2}$. \diamond

4.1.1 Sub-Gaussian random variables

Gaussian random variables are convenient for their nice analytical properties, but a broader class of random variables with similar moment generating functions are known as *sub-Gaussian* random variables.

Definition 4.1. A random variable X is sub-Gaussian with parameter σ^2 if

$$\mathbb{E}[\exp(\lambda(X - \mathbb{E}[X]))] \leq \exp\left(\frac{\lambda^2 \sigma^2}{2}\right)$$

for all $\lambda \in \mathbb{R}$. We also say such a random variable is σ^2 -sub-Gaussian.

Of course, Gaussian random variables satisfy Definition 4.1 with equality. This would be uninteresting if only Gaussian random variables satisfied this property; happily, that is not the case, and we detail several examples.

Example 4.1.5 (Random signs (Rademacher variables)): The random variable X taking values $\{-1, 1\}$ with equal probability is 1-sub-Gaussian. Indeed, we have

$$\mathbb{E}[\exp(\lambda X)] = \frac{1}{2}e^{\lambda} + \frac{1}{2}e^{-\lambda} = \frac{1}{2} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} + \frac{1}{2} \sum_{k=0}^{\infty} \frac{(-\lambda)^k}{k!} = \sum_{k=0}^{\infty} \frac{\lambda^{2k}}{(2k)!} \leq \sum_{k=0}^{\infty} \frac{(\lambda^2)^k}{2^k k!} = \exp\left(\frac{\lambda^2}{2}\right),$$

as claimed. \diamond

Bounded random variables are also sub-Gaussian; indeed, we have the following example.

Example 4.1.6 (Bounded random variables): Suppose that X is bounded, say $X \in [a, b]$. Then Hoeffding's lemma states that

$$\mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] \leq \exp\left(\frac{\lambda^2(b-a)^2}{8}\right),$$

so that X is $(b-a)^2/4$ -sub-Gaussian.

We prove a somewhat weaker statement with a simpler argument, while Exercise 4.1 gives one approach to proving the above statement. First, let $\varepsilon \in \{-1, 1\}$ be a Rademacher variable, so that $\mathbb{P}(\varepsilon = 1) = \mathbb{P}(\varepsilon = -1) = \frac{1}{2}$. We apply a so-called *symmetrization* technique—a common technique in probability theory, statistics, concentration inequalities, and Banach space research—to give a simpler bound. Indeed, let X' be an independent copy of X , so that $\mathbb{E}[X'] = \mathbb{E}[X]$. We have

$$\begin{aligned} \varphi_{X-\mathbb{E}[X]}(\lambda) &= \mathbb{E}[\exp(\lambda(X - \mathbb{E}[X']))] \leq \mathbb{E}[\exp(\lambda(X - X'))] \\ &= \mathbb{E}[\exp(\lambda\varepsilon(X - X'))], \end{aligned}$$

where the inequality follows from Jensen's inequality and the last equality is a consequence of the fact that $X - X'$ is symmetric about 0. Using the result of Example 4.1.5,

$$\mathbb{E}[\exp(\lambda \varepsilon(X - X'))] \leq \mathbb{E}\left[\exp\left(\frac{\lambda^2(X - X')^2}{2}\right)\right] \leq \exp\left(\frac{\lambda^2(b - a)^2}{2}\right),$$

where the final inequality is immediate from the fact that $|X - X'| \leq b - a$. \diamond

While Example 4.1.6 shows how a symmetrization technique can give sub-Gaussian behavior, more sophisticated techniques involving explicitly bounding the logarithm of the moment generating function of X , often by calculations involving *exponential tilts* of its density. In particular, letting X be mean zero for simplicity, if we let

$$\psi(\lambda) = \log \varphi_X(\lambda) = \log \mathbb{E}[e^{\lambda X}],$$

then

$$\psi'(\lambda) = \frac{\mathbb{E}[X e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} \quad \text{and} \quad \psi''(\lambda) = \frac{\mathbb{E}[X^2 e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} - \frac{\mathbb{E}[X e^{\lambda X}]^2}{\mathbb{E}[e^{\lambda X}]^2},$$

where we can interchange the order of taking expectations and derivatives whenever $\psi(\lambda)$ is finite. Notably, if X has density p_X (with respect to any base measure) then the random variable Y_λ with density

$$p_\lambda(y) = \frac{e^{\lambda y}}{\mathbb{E}[e^{\lambda X}]} p_X(y)$$

(with respect to the same base measure) satisfies

$$\psi'(\lambda) = \mathbb{E}[Y_\lambda] \quad \text{and} \quad \psi''(\lambda) = \mathbb{E}[Y_\lambda^2] - \mathbb{E}[Y_\lambda]^2 = \text{Var}(Y_\lambda).$$

One can exploit this in many ways, which the exercises and coming chapters do. As a particular example, we can give sharper sub-Gaussian constants for Bernoulli random variables.

Example 4.1.7 (Bernoulli random variables): Let X be Bernoulli(p), so that $X = 1$ with probability p and $X = 0$ otherwise. Then a strengthening of Hoeffding's lemma (also, essentially, due to Hoeffding) is that

$$\log \mathbb{E}[e^{\lambda(X-p)}] \leq \frac{\sigma^2(p)}{2} \lambda^2 \quad \text{for} \quad \sigma^2(p) := \frac{1 - 2p}{2 \log \frac{1-p}{p}}.$$

Here we take the limits as $p \rightarrow \{0, \frac{1}{2}, 1\}$ and have $\sigma^2(0) = 0$, $\sigma^2(1) = 0$, and $\sigma^2(\frac{1}{2}) = \frac{1}{4}$. Because $p \mapsto \sigma^2(p)$ is concave and symmetric about $p = \frac{1}{2}$, this inequality is always sharper than that of Example 4.1.6. Exercise 4.12 gives one proof of this bound exploiting exponential tilting. \diamond

Chernoff bounds for sub-Gaussian random variables are immediate; indeed, they have the same concentration properties as Gaussian random variables, a consequence of the nice analytical properties of their moment generating functions (that their logarithms are at most quadratic). Thus, using the technique of Example 4.1.4, we obtain the following proposition.

Proposition 4.1.8. *Let X be a σ^2 -sub-Gaussian. Then for all $t \geq 0$ we have*

$$\mathbb{P}(X - \mathbb{E}[X] \geq t) \vee \mathbb{P}(X - \mathbb{E}[X] \leq -t) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right).$$

Chernoff bounds extend naturally to sums of independent random variables, because moment generating functions of sums of independent random variables become products of moment generating functions.

Proposition 4.1.9. *Let X_1, X_2, \dots, X_n be independent σ_i^2 -sub-Gaussian random variables. Then*

$$\mathbb{E} \left[\exp \left(\lambda \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right) \right] \leq \exp \left(\frac{\lambda^2 \sum_{i=1}^n \sigma_i^2}{2} \right) \quad \text{for all } \lambda \in \mathbb{R},$$

that is, $\sum_{i=1}^n X_i$ is $\sum_{i=1}^n \sigma_i^2$ -sub-Gaussian.

Proof We assume w.l.o.g. that the X_i are mean zero. We have by independence that and sub-Gaussianity that

$$\mathbb{E} \left[\exp \left(\lambda \sum_{i=1}^n X_i \right) \right] = \mathbb{E} \left[\exp \left(\lambda \sum_{i=1}^{n-1} X_i \right) \right] \mathbb{E}[\exp(\lambda X_n)] \leq \exp \left(\frac{\lambda^2 \sigma_n^2}{2} \right) \mathbb{E} \left[\exp \left(\lambda \sum_{i=1}^{n-1} X_i \right) \right].$$

Applying this technique inductively to X_{n-1}, \dots, X_1 , we obtain the desired result. \square

Two immediate corollary to Propositions 4.1.8 and 4.1.9 show that sums of sub-Gaussian random variables concentrate around their expectations. We begin with a general concentration inequality.

Corollary 4.1.10. *Let X_i be independent σ_i^2 -sub-Gaussian random variables. Then for all $t \geq 0$*

$$\max \left\{ \mathbb{P} \left(\sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \geq t \right), \mathbb{P} \left(\sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \leq -t \right) \right\} \leq \exp \left(-\frac{t^2}{2 \sum_{i=1}^n \sigma_i^2} \right).$$

Additionally, the classical Hoeffding bound, follows when we couple Example 4.1.6 with Corollary 4.1.10: if $X_i \in [a_i, b_i]$, then

$$\mathbb{P} \left(\sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \geq t \right) \leq \exp \left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2} \right).$$

To give another interpretation of these inequalities, let us assume that X_i are independent and σ^2 -sub-Gaussian. Then we have that

$$\mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \geq t \right) \leq \exp \left(-\frac{nt^2}{2\sigma^2} \right),$$

or, for $\delta \in (0, 1)$, setting $\exp(-\frac{nt^2}{2\sigma^2}) = \delta$ or $t = \frac{\sqrt{2\sigma^2 \log \frac{1}{\delta}}}{\sqrt{n}}$, we have that

$$\frac{1}{n} \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \leq \frac{\sqrt{2\sigma^2 \log \frac{1}{\delta}}}{\sqrt{n}} \quad \text{with probability at least } 1 - \delta.$$

There are a variety of other conditions equivalent to sub-Gaussianity, which we capture in the following theorem.

Theorem 4.1.11. *Let X be a random variable and $\sigma^2 \geq 0$. The following statements are all equivalent, meaning that there are numerical constant factors K_j such that if one statement (i) holds with parameter K_i , then statement (j) holds with parameter $K_j \leq CK_i$, where C is a numerical constant.*

(1) *Sub-gaussian tails:* $\mathbb{P}(|X| \geq t) \leq 2 \exp(-\frac{t^2}{K_1 \sigma^2})$ for all $t \geq 0$.

(2) *Sub-gaussian moments:* $\mathbb{E}[|X|^k]^{1/k} \leq K_2 \sigma \sqrt{k}$ for all k .

(3) *Super-exponential moment:* $\mathbb{E}[\exp(X^2/(K_3 \sigma^2))] \leq e$.

If in addition X is mean zero, each of these is equivalent to

(4) *Sub-gaussian moment generating function:* $\mathbb{E}[\exp(\lambda X)] \leq \exp(K_4 \lambda^2 \sigma^2)$ for all $\lambda \in \mathbb{R}$.

Particularly, (1) implies (2) with $K_1 = 1$ and $K_2 \leq e^{1/e}$; (2) implies (3) with $K_2 = 1$ and $K_3 = e \sqrt{\frac{2}{e-1}} < 3$; (3) implies (1) with $K_3 = 1$ and $K_1 = 1/\log 2$. For the last part, (3) implies (4) with $K_3 = 1$ and $K_4 \leq \frac{3}{4}$, while (4) implies (1) with $K_4 = \frac{1}{2}$ and $K_1 \leq 2$.

This result is standard in the literature on concentration and random variables; see Section 4.4.1 for a proof.

For completeness, we can give a tighter result than part (3) of the preceding theorem, giving a concrete upper bound on squares of sub-Gaussian random variables. The technique used in the example, to introduce an independent random variable for auxiliary randomization, is a common and useful technique in probabilistic arguments (similar to our use of symmetrization in Example 4.1.6).

Example 4.1.12 (Sub-Gaussian squares): Let X be a mean-zero σ^2 -sub-Gaussian random variable. Then

$$\mathbb{E}[\exp(\lambda X^2)] \leq \frac{1}{[1 - 2\sigma^2 \lambda]_+^{\frac{1}{2}}}, \quad (4.1.2)$$

and expression (4.1.2) holds with equality for $X \sim \mathcal{N}(0, \sigma^2)$.

To see this result, we focus on the Gaussian case first and assume (for this case) without loss of generality (by scaling) that $\sigma^2 = 1$. Assuming that $\lambda < \frac{1}{2}$, we have

$$\mathbb{E}[\exp(\lambda Z^2)] = \int \frac{1}{\sqrt{2\pi}} e^{-(\frac{1}{2}-\lambda)z^2} dz = \int \frac{1}{\sqrt{2\pi}} e^{-\frac{1-2\lambda}{2}z^2} dz = \frac{\sqrt{2\pi}}{\sqrt{1-2\lambda}} \frac{1}{\sqrt{2\pi}},$$

the final equality a consequence of the fact that (as we know for normal random variables) $\int e^{-\frac{1}{2\sigma^2}z^2} dz = \sqrt{2\pi\sigma^2}$. When $\lambda \geq \frac{1}{2}$, the above integrals are all infinite, giving the equality in expression (4.1.2).

For the more general inequality, we recall that if Z is an independent $\mathcal{N}(0, 1)$ random variable, then $\mathbb{E}[\exp(tZ)] = \exp(\frac{t^2}{2})$, and so

$$\mathbb{E}[\exp(\lambda X^2)] = \mathbb{E}[\exp(\sqrt{2\lambda} X Z)] \stackrel{(i)}{\leq} \mathbb{E}[\exp(\lambda \sigma^2 Z^2)] \stackrel{(ii)}{=} \frac{1}{[1 - 2\sigma^2 \lambda]_+^{\frac{1}{2}}},$$

where inequality (i) follows because X is sub-Gaussian, and (ii) because $Z \sim \mathcal{N}(0, 1)$. \diamond

4.1.2 Sub-exponential random variables

A slightly weaker condition than sub-Gaussianity is for a random variable to be *sub-exponential*, which—for a mean-zero random variable—means that its moment generating function exists in a neighborhood of zero.

Definition 4.2. A random variable X is sub-exponential with parameters (τ^2, b) if for all λ such that $|\lambda| \leq 1/b$,

$$\mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] \leq \exp\left(\frac{\lambda^2 \tau^2}{2}\right).$$

It is clear from Definition 4.2 that a σ^2 -sub-Gaussian random variable is $(\sigma^2, 0)$ -sub-exponential.

A variety of random variables are sub-exponential. As a first example, χ^2 -random variables are sub-exponential with constant values for τ and b :

Example 4.1.13: Let $X = Z^2$, where $Z \sim \mathcal{N}(0, 1)$. We claim that

$$\mathbb{E}[\exp(\lambda(X - \mathbb{E}[X]))] \leq \exp(2\lambda^2) \quad \text{for } \lambda \leq \frac{1}{4}. \quad (4.1.3)$$

Indeed, for $\lambda < \frac{1}{2}$ we have (recall Example 4.1.12) that

$$\mathbb{E}[\exp(\lambda(Z^2 - \mathbb{E}[Z^2]))] = \exp\left(-\frac{1}{2} \log(1 - 2\lambda) - \lambda\right) \stackrel{(*)}{\leq} \exp(\lambda + 2\lambda^2 - \lambda)$$

where inequality $(*)$ holds for $\lambda \leq \frac{1}{4}$, because $-\log(1 - 2\lambda) \leq 2\lambda + 4\lambda^2$ for $\lambda \leq \frac{1}{4}$. \diamond

As a second example, we can show that bounded random variables are sub-exponential. It is clear that this is the case as they are also sub-Gaussian; however, in many cases, it is possible to show that their parameters yield much tighter control over deviations than is possible using only sub-Gaussian techniques.

Example 4.1.14 (Bounded random variables are sub-exponential): Suppose that X is a mean zero random variable taking values in $[-b, b]$ with variance $\sigma^2 = \mathbb{E}[X^2]$ (note that we are guaranteed that $\sigma^2 \leq b^2$ in this case). We claim that

$$\mathbb{E}[\exp(\lambda X)] \leq \exp\left(\frac{3\lambda^2 \sigma^2}{5}\right) \quad \text{for } |\lambda| \leq \frac{1}{2b}. \quad (4.1.4)$$

To see this, we expand e^z via

$$e^z = 1 + z + \frac{z^2}{2} \sum_{k=2}^{\infty} \frac{2z^{k-2}}{k!} = 1 + z + \frac{z^2}{2} \sum_{k=0}^{\infty} \frac{2}{(k+2)!} z^k.$$

For $k \geq 0$, we have $\frac{2}{(k+2)!} \leq \frac{1}{3^k}$, so that $\sum_{k=0}^{\infty} \frac{2}{(k+2)!} z^k \leq \sum_{k=0}^{\infty} |z/3|^k = [1 - |z|/3]_+^{-1}$. Thus

$$e^z \leq 1 + z + \frac{1}{[1 - |z|/3]_+} \frac{z^2}{2},$$

and as $|X| \leq b$ and $|\lambda| < \frac{3}{b}$ we therefore obtain

$$\mathbb{E}[\exp(\lambda X)] \leq 1 + \mathbb{E}[\lambda X] + \mathbb{E}\left[\frac{\lambda^2 X^2}{2[1 - |\lambda X|/3]_+}\right] \leq 1 + \frac{1}{1 - |\lambda|b/3} \frac{\lambda^2 \sigma^2}{2}.$$

Letting $|\lambda| \leq \frac{1}{2b}$ implies $\frac{1}{1-|\lambda|b/3} \leq \frac{6}{5}$, and using that $1+x \leq e^x$ gives the result.

It is possible to give a slightly tighter result for $\lambda \geq 0$. In this case, we have the bound

$$\mathbb{E}[\exp(\lambda X)] \leq 1 + \frac{\lambda^2 \sigma^2}{2} + \lambda^2 \sigma^2 \sum_{k=3}^{\infty} \frac{\lambda^{k-2} b^{k-2}}{k!} = 1 + \frac{\sigma^2}{b^2} (e^{\lambda b} - 1 - \lambda b).$$

Then using that $1+x \leq e^x$, we obtain *Bennett's moment generating inequality*, which is that

$$\mathbb{E}[e^{\lambda X}] \leq \exp\left(\frac{\sigma^2}{b^2} (e^{\lambda b} - 1 - \lambda b)\right) \quad \text{for } \lambda \geq 0. \quad (4.1.5)$$

Inequality (4.1.5) always holds, and for λb near 0, we have $e^{\lambda b} - 1 - \lambda b \approx \frac{\lambda^2 b^2}{2}$. \diamond

In particular, if the variance $\sigma^2 \ll b^2$, the absolute bound on X , inequality (4.1.4) gives much tighter control on the moment generating function of X than typical sub-Gaussian bounds based only on the fact that $X \in [-b, b]$ allow.

More broadly, we can show a result similar to Theorem 4.1.11.

Theorem 4.1.15. *Let X be a random variable and $\sigma \geq 0$. Then—in the sense of Theorem 4.1.11—the following statements are all equivalent for suitable numerical constants K_1, \dots, K_4 .*

- (1) *Sub-exponential tails:* $\mathbb{P}(|X| \geq t) \leq 2 \exp(-\frac{t}{K_1 \sigma})$ for all $t \geq 0$
- (2) *Sub-exponential moments:* $\mathbb{E}[|X|^k]^{1/k} \leq K_2 \sigma k$ for all $k \geq 1$.
- (3) *Existence of moment generating function:* $\mathbb{E}[\exp(X/(K_3 \sigma))] \leq e$ and $\mathbb{E}[\exp(-X/(K_3 \sigma))] \leq e$.

If in addition X is mean zero, each of these is equivalent to

- (4) *Sub-exponential moment generating function:* $\mathbb{E}[\exp(\lambda X)] \leq \exp(K_4 \lambda^2 \sigma^2)$ for $|\lambda| \leq K'_4 / \sigma$.

In particular, if (2) holds with $K_2 = 1$, then (4) holds with $K_4 = 2e^2$ and $K'_4 = \frac{1}{2e}$.

See Section 4.4.2 for the proof, which is similar to that for Theorem 4.1.11.

While the concentration properties of sub-exponential random variables are not quite so nice as those for sub-Gaussian random variables (recall Hoeffding's inequality, Corollary 4.1.10), we can give sharp tail bounds for sub-exponential random variables. We first give a simple bound on deviation probabilities.

Proposition 4.1.16. *Let X be a mean-zero (τ^2, b) -sub-exponential random variable. Then for all $t \geq 0$,*

$$\mathbb{P}(X \geq t) \vee \mathbb{P}(X \leq -t) \leq \exp\left(-\frac{1}{2} \min\left\{\frac{t^2}{\tau^2}, \frac{t}{b}\right\}\right).$$

Proof The proof is an application of the Chernoff bound technique; we prove only the upper tail as the lower tail is similar. We have

$$\mathbb{P}(X \geq t) \leq \frac{\mathbb{E}[e^{\lambda X}]}{e^{\lambda t}} \stackrel{(i)}{\leq} \exp\left(\frac{\lambda^2 \tau^2}{2} - \lambda t\right),$$

inequality (i) holding for $|\lambda| \leq 1/b$. To minimize the last term in λ , we take $\lambda = \min\{\frac{t}{\tau^2}, 1/b\}$, which gives the result. \square

Comparing with sub-Gaussian random variables, which have $b = 0$, we see that Proposition 4.1.16 gives a similar result for small t —essentially the same concentration sub-Gaussian random variables—while for large t , the tails decrease only exponentially in t .

We can also give a tensorization identity similar to Proposition 4.1.9.

Proposition 4.1.17. *Let X_1, \dots, X_n be independent mean-zero sub-exponential random variables, where X_i is (σ_i^2, b_i) -sub-exponential. Then for any vector $a \in \mathbb{R}^n$, we have*

$$\mathbb{E} \left[\exp \left(\lambda \sum_{i=1}^n a_i X_i \right) \right] \leq \exp \left(\frac{\lambda^2 \sum_{i=1}^n a_i^2 \sigma_i^2}{2} \right) \quad \text{for } |\lambda| \leq \frac{1}{b_*},$$

where $b_* = \max_i b_i |a_i|$. That is, $\langle a, X \rangle$ is $(\sum_{i=1}^n a_i^2 \sigma_i^2, \max_i b_i |a_i|)$ -sub-exponential.

Proof We apply an inductive technique similar to that used in the proof of Proposition 4.1.9. First, for any fixed i , we know that if $|\lambda| \leq \frac{1}{b_i |a_i|}$, then $|a_i \lambda| \leq \frac{1}{b_i}$ and so

$$\mathbb{E}[\exp(\lambda a_i X_i)] \leq \exp \left(\frac{\lambda^2 a_i^2 \sigma_i^2}{2} \right).$$

Now, we inductively apply the preceding inequality, which applies so long as $|\lambda| \leq \frac{1}{b_i |a_i|}$ for all i . We have

$$\mathbb{E} \left[\exp \left(\lambda \sum_{i=1}^n a_i X_i \right) \right] = \prod_{i=1}^n \mathbb{E}[\exp(\lambda a_i X_i)] \leq \prod_{i=1}^n \exp \left(\frac{\lambda^2 a_i^2 \sigma_i^2}{2} \right),$$

which is our desired result. \square

As in the case of sub-Gaussian random variables, a combination of the tensorization property—that the moment generating functions of sums of sub-exponential random variables are well-behaved—of Proposition 4.1.17 and the concentration inequality (4.1.16) immediately yields the following Bernstein-type inequality. (See also Vershynin [187].)

Corollary 4.1.18. *Let X_1, \dots, X_n be independent mean-zero (σ_i^2, b_i) -sub-exponential random variables (Definition 4.2). Define $b_* := \max_i b_i$. Then for all $t \geq 0$ and all vectors $a \in \mathbb{R}^n$, we have*

$$\mathbb{P} \left(\sum_{i=1}^n a_i X_i \geq t \right) \vee \mathbb{P} \left(\sum_{i=1}^n a_i X_i \leq -t \right) \leq \exp \left(-\frac{1}{2} \min \left\{ \frac{t^2}{\sum_{i=1}^n a_i^2 \sigma_i^2}, \frac{t}{b_* \|a\|_\infty} \right\} \right).$$

It is instructive to study the structure of the bound of Corollary 4.1.18. Notably, the bound is similar to the Hoeffding-type bound of Corollary 4.1.10 (holding for σ^2 -sub-Gaussian random variables) that

$$\mathbb{P} \left(\sum_{i=1}^n a_i X_i \geq t \right) \leq \exp \left(-\frac{t^2}{2 \|a\|_2^2 \sigma^2} \right),$$

so that for small t , Corollary 4.1.18 gives sub-Gaussian tail behavior. For large t , the bound is weaker. However, in many cases, Corollary 4.1.18 can give finer control than naive sub-Gaussian bounds. Indeed, suppose that the random variables X_i are i.i.d., mean zero, and satisfy $X_i \in [-b, b]$ with probability 1, but have variance $\sigma^2 = \mathbb{E}[X_i^2] \leq b^2$ as in Example 4.1.14. Then Corollary 4.1.18 implies that

$$\mathbb{P} \left(\sum_{i=1}^n a_i X_i \geq t \right) \leq \exp \left(-\frac{1}{2} \min \left\{ \frac{5}{6} \frac{t^2}{\sigma^2 \|a\|_2^2}, \frac{t}{2b \|a\|_\infty} \right\} \right). \quad (4.1.6)$$

When applied to a standard mean (and with a minor simplification that $5/12 < 1/3$) with $a_i = \frac{1}{n}$, we obtain the bound that $\frac{1}{n} \sum_{i=1}^n X_i \leq t$ with probability at least $1 - \exp(-n \min\{\frac{t^2}{3\sigma^2}, \frac{t}{4b}\})$. Written differently, we take $t = \max\{\sigma\sqrt{\frac{3 \log \frac{1}{\delta}}{n}}, \frac{4b \log \frac{1}{\delta}}{n}\}$ to obtain

$$\frac{1}{n} \sum_{i=1}^n X_i \leq \max \left\{ \sigma \frac{\sqrt{3 \log \frac{1}{\delta}}}{\sqrt{n}}, \frac{4b \log \frac{1}{\delta}}{n} \right\} \quad \text{with probability } 1 - \delta.$$

The sharpest such bound possible via more naive Hoeffding-type bounds is $b\sqrt{2 \log \frac{1}{\delta}}/\sqrt{n}$, which has substantially worse scaling.

The exercises ask you to work out further variants of these results, including the sub-exponential behavior of quadratic forms of Gaussian random vectors. As one particular example, Exercises 4.10 and 4.11 work through the details of proving the following corollary.

Corollary 4.1.19. *Let $Z \sim \mathcal{N}(0, 1)$. Then for any $\mu \in \mathbb{R}$, $(\mu + Z)^2$ is $(4(1+2\mu^2), 4)$ -sub-exponential, and more precisely,*

$$\mathbb{E} \left[\exp \left(\lambda \left((\mu + Z)^2 - (\mu^2 + 1) \right) \right) \right] \leq \exp \left(\frac{2\lambda^2 \mu^2}{1 - 2\lambda} + \frac{\lambda^2}{[1 - 2|\lambda|]_+} \right).$$

Additionally, if $Z \sim \mathcal{N}(0, I)$, then for any matrix A and vector b , $\|AZ - b\|_2^2$ is sub-exponential with

$$\mathbb{E} \left[\exp \left(\lambda \left(\|AZ - b\|_2^2 - \|A\|_{\text{Fr}}^2 - \|b\|_2^2 \right) \right) \right] \leq \exp \left(2\lambda^2 (\|A\|_{\text{Fr}}^2 + 2\|b\|_2^2) \right) \quad \text{for } |\lambda| \leq \frac{1}{4\|A\|_{\text{op}}^2}.$$

Further conditions and examples

There are a number of examples and conditions sufficient for random variables to be sub-exponential. One common condition, the so-called *Bernstein* condition, controls the higher moments of a random variable X by its variance. In this case, we say that X satisfies the b -Bernstein condition if

$$|\mathbb{E}[(X - \mu)^k]| \leq \frac{k!}{2} \sigma^2 b^{k-2} \quad \text{for } k = 3, 4, \dots, \quad (4.1.7)$$

where $\mu = \mathbb{E}[X]$ and $\sigma^2 = \text{Var}(X) = \mathbb{E}[X^2] - \mu^2$. In this case, the following lemma controls the moment generating function of X . This result is essentially present in Theorem 4.1.15, but it provides somewhat tighter control with precise constants.

Lemma 4.1.20. *Let X be a random variable satisfying the Bernstein condition (4.1.7). Then*

$$\mathbb{E} \left[e^{\lambda(X-\mu)} \right] \leq \exp \left(\frac{\lambda^2 \sigma^2}{2(1 - b|\lambda|)} \right) \quad \text{for } |\lambda| \leq \frac{1}{b}.$$

Said differently, a random variable satisfying Condition (4.1.7) is $(\sqrt{2}\sigma, b/2)$ -sub-exponential.

Proof Without loss of generality we assume $\mu = 0$. We expand the moment generating function by noting that

$$\begin{aligned} \mathbb{E}[e^{\lambda X}] &= 1 + \frac{\lambda^2 \sigma^2}{2} + \sum_{k=3}^{\infty} \frac{\lambda^k \mathbb{E}[X^k]}{k!} \stackrel{(i)}{\leq} 1 + \frac{\lambda^2 \sigma^2}{2} + \frac{\lambda^2 \sigma^2}{2} \sum_{k=3}^{\infty} |\lambda b|^{k-2} \\ &= 1 + \frac{\lambda^2 \sigma^2}{2} \frac{1}{[1 - b|\lambda|]_+} \end{aligned}$$

where inequality (i) used the Bernstein condition (4.1.7). Noting that $1+x \leq e^x$ gives the result. \square

As one final example, we return to Bennett's inequality (4.1.5) from Example 4.1.14.

Proposition 4.1.21 (Bennett's inequality). *Let X_i be independent mean-zero random variables with $\text{Var}(X_i) = \sigma_i^2$ and $|X_i| \leq b$. Then for $h(t) := (1+t) \log(1+t) - t$ and $\sigma^2 := \sum_{i=1}^n \sigma_i^2$, we have*

$$\mathbb{P} \left(\sum_{i=1}^n X_i \geq t \right) \leq \exp \left(-\frac{\sigma^2}{b^2} h \left(\frac{bt}{\sigma^2} \right) \right).$$

Proof We assume without loss of generality that $\mathbb{E}[X] = 0$. Using the standard Chernoff bound argument coupled with inequality (4.1.5), we see that

$$\mathbb{P} \left(\sum_{i=1}^n X_i \geq t \sum \right) \leq \exp \left(\sum_{i=1}^n \frac{\sigma_i^2}{b^2} (e^{\lambda b} - 1 - \lambda b) - \lambda t \right).$$

Letting $h(t) = (1+t) \log(1+t) - t$ as in the statement of the proposition and $\sigma^2 = \sum_{i=1}^n \sigma_i^2$, we minimize over $\lambda \geq 0$, setting $\lambda = \frac{1}{b} \log(1 + \frac{bt}{\sigma^2})$. Substituting into our Chernoff bound application gives the proposition. \square

A slightly more intuitive writing of Bennett's inequality is to use averages, in which case for $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \sigma_i^2$ the average of the variances,

$$\mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n X_i \geq t \right) \leq \exp \left(-\frac{n\sigma^2}{b} h \left(\frac{bt}{\sigma^2} \right) \right).$$

It is possible to show that

$$\frac{n\sigma^2}{b} h \left(\frac{bt}{\sigma^2} \right) \geq \frac{nt^2}{2\sigma^2 + \frac{2}{3}bt},$$

which gives rise to the classical Bernstein inequality that

$$\mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n X_i \geq t \right) \leq \exp \left(-\frac{nt^2}{2\sigma^2 + \frac{2}{3}bt} \right). \quad (4.1.8)$$

4.1.3 Orlicz norms

Sub-Gaussian and sub-exponential random variables are examples of a broader class of random variables belonging to what are known as *Orlicz-spaces*. For these, we take any convex function $\psi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ with $\psi(0) = 0$ and $\psi(t) \rightarrow \infty$ as $t \uparrow \infty$, a class called the Orlicz functions. Then the Orlicz norm of a random variable X is

$$\|X\|_\psi := \inf \{t > 0 \mid \mathbb{E}[\psi(|X|/t)] \leq 1\}. \quad (4.1.9)$$

That this is a norm is not completely trivial, though a few properties are immediate: clearly $\|aX\|_\psi = |a| \|X\|_\psi$, and we have $\|X\|_\psi = 0$ if and only if $X = 0$ with probability 1. The key result is that in fact, $\|\cdot\|_\psi$ is actually convex, which then guarantees that it is a norm.

Proposition 4.1.22. *The function $\|\cdot\|_\psi$ is convex on the space of random variables.*

Proof Because ψ is convex and non-decreasing, $x \mapsto \psi(|x|)$ is convex as well. (Convince yourself of this.) Thus, its *perspective transform* $\text{pers}(\psi)(t, |x|) := t\psi(|x|/t)$ is jointly convex in both $t \geq 0$ and x (see Appendix B.3.3). This joint convexity of $\text{pers}(\psi)$ implies that for any random variables X_0 and X_1 and t_0, t_1 ,

$$\mathbb{E}[\text{pers}(\psi)(\lambda t_0 + (1 - \lambda)t_1, |\lambda X_0 + (1 - \lambda)X_1|)] \leq \lambda \mathbb{E}[\text{pers}(\psi)(t_0, |X_0|)] + (1 - \lambda) \mathbb{E}[\text{pers}(\psi)(t_1, |X_1|)].$$

Now note that $\mathbb{E}[\psi(|X|/t)] \leq 1$ if and only if $t\mathbb{E}[\psi(|X|/t)] \leq t$. \square

Because $\|\cdot\|_\psi$ is convex and positively homogeneous, we certainly have

$$\|X + Y\|_\psi = 2 \|(X + Y)/2\|_\psi \leq \|X\|_\psi + \|Y\|_\psi,$$

that is, the triangle inequality holds. This implies that centering a variable can never increase its norm by much:

$$\|X - \mathbb{E}[X]\|_\psi \leq \|X\|_\psi + \|\mathbb{E}[X]\|_\psi \leq \|X\|_\psi + \|X\|_\psi$$

by Jensen's inequality, so that $\|X - \mathbb{E}[X]\|_\psi \leq 2\|X\|_\psi$.

We can recover several standard norms on random variables, including some we have already implicitly used. The first are the classical L^p norms, where we take $\psi(t) = t^p$, where we see that

$$\inf\{t > 0 \mid \mathbb{E}[|X|^p/t^p] \leq 1\} = \mathbb{E}[|X|^p]^{1/p}.$$

We also have what we term the *sub-Gaussian* and *sub-Exponential* norms, which we denote by considering the functions

$$\psi_p(x) := \exp(|x|^p) - 1.$$

These induce the *Orlicz ψ_p -norms*, as for $p \geq 1$, these are convex (as they are the composition of the increasing convex function $\exp(\cdot)$ applied to the nonnegative convex function $|\cdot|^p$). Theorem 4.1.11 shows that we have a *sub-Gaussian* norm

$$\|X\|_{\psi_2} := \inf\{t > 0 \mid \mathbb{E}[\exp(X^2/t^2)] \leq 2\}, \quad (4.1.10)$$

while Theorem 4.1.15 shows a *sub-exponential* norm (or Orlicz ψ_1 -norm)

$$\|X\|_{\psi_1} := \inf\{t > 0 \mid \mathbb{E}[\exp(|X|/t)] \leq 2\}. \quad (4.1.11)$$

Many relationships follow immediately from the definitions (4.1.10) and (4.1.11). For example, the definition of the ψ_p -norms immediately implies that a sub-Gaussian random variable (whether or not it is mean zero) has a sub-exponential square:

Lemma 4.1.23. *A random variable X is sub-Gaussian if and only if X^2 is sub-exponential, as*

$$\|X\|_{\psi_2}^2 = \|X^2\|_{\psi_1}.$$

Additionally,

$$\|XY\|_{\psi_1} \leq \|X\|_{\psi_2} \|Y\|_{\psi_2}.$$

Proof We prove only the second statement. Because $xy \leq \frac{x^2}{2\eta} + \frac{\eta y^2}{2}$ for any x, y , and any $\eta > 0$, for any $t > 0$ we have

$$\mathbb{E}[\exp(|XY|/t)] \leq \mathbb{E} \left[\exp \left(\frac{X^2}{2\eta t} + \frac{\eta Y^2}{2t} \right) \right] \leq \mathbb{E}[\exp(X^2/\eta t)]^{1/2} \mathbb{E}[\exp(\eta Y^2/t)]^{1/2}$$

by the Cauchy-Schwarz inequality. In particular, if we take $t = \|X\|_{\psi_2} \|Y\|_{\psi_2}$, then the choice $\eta = \|X\|_{\psi_2}^2 / \|Y\|_{\psi_2}^2$ gives $\mathbb{E}[\exp(X^2/\eta t)] \leq 2$ and $\mathbb{E}[\exp(\eta Y^2/t)] \leq 2$, so that $\mathbb{E}[\exp(|XY|/t)] \leq 2$. \square

By tracing through the arguments in the proofs of Theorems 4.1.11 and 4.1.15, we can also see that we have the equivalences

$$\|X\|_{\psi_2} \asymp \sup_{k \in \mathbb{N}} \frac{1}{\sqrt{k}} \mathbb{E}[|X|^k]^{1/k} \quad \text{and} \quad \|X\|_{\psi_1} \asymp \sup_{k \in \mathbb{N}} \frac{1}{k} \mathbb{E}[|X|^k]^{1/k},$$

where \asymp denotes upper and lower bounds by numerical constants.

The arguments we use to prove Theorems 4.1.11 and 4.1.15 also show the following result, which gives explicit constants connecting sub-exponential behavior with the ψ_1 -norm.

Corollary 4.1.24. *Let X be any random variable with $\|X\|_{\psi_1} < \infty$. Then for all $t \geq 0$,*

$$\mathbb{P}(|X| \geq t) \leq 2 \exp(-t / \|X\|_{\psi_1})$$

and if $\mathbb{E}[X] = 0$, then X is $(8\|X\|_{\psi_1}^2, 2\|X\|_{\psi_1})$ -sub-exponential.

Proof The first statement is nearly trivial: we have by the Chernoff bounding method that

$$\mathbb{P}(|X| \geq t) \leq \mathbb{E} \left[\exp(|X| / \|X\|_{\psi_1}) \right] \exp(-t / \|X\|_{\psi_1}) \leq 2 \exp(-t / \|X\|_{\psi_1})$$

by definition of the ψ_1 -norm. For the second, we mimic the proof of Theorem 4.1.15: because $\mathbb{E}[Z] = \int_0^\infty \mathbb{P}(Z \geq t) dt$ for $Z \geq 0$, we have

$$\frac{\mathbb{E}[|X|^k]}{\|X\|_{\psi_1}^k} \leq \int_0^\infty \mathbb{P}(|X| / \|X\|_{\psi_1} \geq t^{1/k}) dt = k \int_0^\infty \mathbb{P}(|X| / \|X\|_{\psi_1} \geq u) u^{k-1} du \leq 2k \int_0^\infty u^{k-1} e^{-u} du$$

using the substitution $u^k = t$. Rearranging yields $\mathbb{E}[|X|^k] \leq 2\|X\|_{\psi_1}^k \Gamma(k+1) = 2\|X\|_{\psi_1}^k k!$. Then computing the moment generating function, we obtain

$$\mathbb{E} \left[\exp(\lambda X / \|X\|_{\psi_1}) \right] \leq 1 + \sum_{k=2}^\infty \frac{\lambda^k \mathbb{E}[|X|^k]}{\|X\|_{\psi_1}^k k!} \leq 1 + 2 \sum_{k=2}^\infty \lambda^k = 1 + \frac{2\lambda^2}{1-|\lambda|}$$

for $|\lambda| < 1$. For $|\lambda| \leq \frac{1}{2}$, we use $1+x \leq e^x$ to obtain $\mathbb{E}[\exp(\lambda X / \|X\|_{\psi_1})] \leq \exp(4\lambda^2)$, which is the desired result. \square

4.1.4 First applications of concentration: random projections

In this section, we investigate the use of concentration inequalities in random projections. As motivation, consider nearest-neighbor (or k -nearest-neighbor) classification schemes. We have a sequence of data points as pairs (u_i, y_i) , where the vectors $u_i \in \mathbb{R}^d$ have labels $y_i \in \{1, \dots, L\}$, where L is the number of possible labels. Given a new point $u \in \mathbb{R}^d$ that we wish to label, we find the k -nearest neighbors to u in the sample $\{(u_i, y_i)\}_{i=1}^n$, then assign u the majority label of these k -nearest neighbors (ties are broken randomly). Unfortunately, it can be prohibitively expensive to store high-dimensional vectors and search over large datasets to find near vectors; this has motivated a line of work in computer science on fast methods for nearest neighbors based on reducing the dimension while preserving essential aspects of the dataset. This line of research begins with Indyk and Motwani [117], and continuing through a variety of other works, including Indyk [116] and work on locality-sensitive hashing by Andoni et al. [7], among others. The original approach is due to Johnson and Lindenstrauss, who used the results in the study of Banach spaces [123]; our proof follows a standard argument.

The most specific variant of this problem is as follows: we have n points u_1, \dots, u_n , and we could like to construct a mapping $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$, where $m \ll d$, such that

$$\|\Phi u_i - \Phi u_j\|^2 \in (1 \pm \epsilon) \|u_i - u_j\|^2.$$

Depending on the norm chosen, this task may be impossible; for the Euclidean (ℓ_2) norm, however, such an embedding is easy to construct using Gaussian random variables and with $m = O(\frac{1}{\epsilon^2} \log n)$. This embedding is known as the Johnson-Lindenstrauss embedding. Note that this size m is *independent* of the dimension d , only depending on the number of points n .

Example 4.1.25 (Johnson-Lindenstrauss): Let the matrix $\Phi \in \mathbb{R}^{m \times d}$ be defined as follows:

$$\Phi_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1/m),$$

and let $\Phi_i \in \mathbb{R}^d$ denote the i th row of this matrix. We claim that

$$m \geq \frac{8}{\epsilon^2} \left[2 \log n + \log \frac{1}{\delta} \right] \quad \text{implies} \quad \|\Phi u_i - \Phi u_j\|_2^2 \in (1 \pm \epsilon) \|u_i - u_j\|_2^2$$

for all pairs u_i, u_j with probability at least $1 - \delta$. In particular, $m \gtrsim \frac{\log n}{\epsilon^2}$ is sufficient to achieve accurate dimension reduction with high probability.

To see this, note that for any fixed vector u ,

$$\frac{\langle \Phi_i, u \rangle}{\|u\|_2} \sim \mathcal{N}(0, 1/m), \quad \text{and} \quad \frac{\|\Phi u\|_2^2}{\|u\|_2^2} = \sum_{i=1}^m \langle \Phi_i, u / \|u\|_2 \rangle^2$$

is a sum of independent scaled χ^2 -random variables. In particular, we have $\mathbb{E}[\|\Phi u / \|u\|_2\|_2^2] = 1$, and using the χ^2 -concentration result of Example 4.1.13 yields

$$\begin{aligned} \mathbb{P} \left(\left| \|\Phi u\|_2^2 / \|u\|_2^2 - 1 \right| \geq \epsilon \right) &= \mathbb{P} \left(m \left| \|\Phi u\|_2^2 / \|u\|_2^2 - 1 \right| \geq m\epsilon \right) \\ &\leq 2 \inf_{|\lambda| \leq \frac{1}{4}} \exp(2m\lambda^2 - \lambda m\epsilon) = 2 \exp \left(-\frac{m\epsilon^2}{8} \right), \end{aligned}$$

the last inequality holding for $\epsilon \in [0, 1]$. Now, using the union bound applied to each of the pairs (u_i, u_j) in the sample, we have

$$\mathbb{P}\left(\text{there exist } i \neq j \text{ s.t. } \left| \|\Phi(u_i - u_j)\|_2^2 - \|u_i - u_j\|_2^2 \right| \geq \epsilon \|u_i - u_j\|_2^2\right) \leq 2 \binom{n}{2} \exp\left(-\frac{m\epsilon^2}{8}\right).$$

Taking $m \geq \frac{8}{\epsilon^2} \log \frac{n^2}{\delta} = \frac{16}{\epsilon^2} \log n + \frac{8}{\epsilon^2} \log \frac{1}{\delta}$ yields that with probability at least $1 - \delta$, we have $\|\Phi u_i - \Phi u_j\|_2^2 \in (1 \pm \epsilon) \|u_i - u_j\|_2^2$. \diamond

Computing low-dimensional embeddings of high-dimensional data is an area of active research, and more recent work has shown how to achieve sharper constants [63] and how to use more structured matrices to allow substantially faster computation of the embeddings Φu (see, for example, Achlioptas [2] for early work in this direction, and Ailon and Chazelle [5] for the so-called “Fast Johnson-Lindenstrauss transform”).

4.1.5 A second application of concentration: codebook generation

We now consider a (very simplified and essentially un-implementable) view of encoding a signal for transmission and generation of a codebook for transmitting said signal. Suppose that we have a set of words, or signals, that we wish to transmit; let us index them by $i \in \{1, \dots, m\}$, so that there are m total signals we wish to communicate across a *binary symmetric channel* Q , meaning that given an input bit $x \in \{0, 1\}$, Q outputs a $z \in \{0, 1\}$ with $Q(Z = x | x) = 1 - \epsilon$ and $Q(Z = 1 - x | x) = \epsilon$, for some $\epsilon < \frac{1}{2}$. (For simplicity, we assume Q is *memoryless*, meaning that when the channel is used multiple times on a sequence x_1, \dots, x_n , its outputs Z_1, \dots, Z_n are conditionally independent: $Q(Z_{1:n} = z_{1:n} | x_{1:n}) = Q(Z_1 = z_1 | x_1) \cdots Q(Z_n = z_n | x_n)$.)

We consider a simplified block coding scheme, where we for each i we associate a codeword $x_i \in \{0, 1\}^d$, where d is a dimension (block length) to be chosen. Upon sending the codeword over the channel, and receiving some $z^{\text{rec}} \in \{0, 1\}^d$, we decode by choosing

$$i^* \in \operatorname{argmax}_{i \in [m]} Q(Z = z^{\text{rec}} | x_i) = \operatorname{argmin}_{i \in [m]} \|z^{\text{rec}} - x_i\|_1, \quad (4.1.12)$$

the maximum likelihood decoder. We now investigate how to choose a collection $\{x_1, \dots, x_m\}$ of such codewords and give finite sample bounds on its probability of error. In fact, by using concentration inequalities, we can show that a randomly drawn codebook of fairly small dimension is likely to enjoy good performance.

Intuitively, if our codebook $\{x_1, \dots, x_m\} \subset \{0, 1\}^d$ is *well-separated*, meaning that each pair of words x_i, x_k satisfies $\|x_i - x_k\|_1 \geq cd$ for some numerical constant $c > 0$, we should be unlikely to make a mistake. Let us make this precise. We mistake word i for word k only if the received signal Z satisfies $\|Z - x_i\|_1 \geq \|Z - x_k\|_1$, and letting $J = \{j \in [d] : x_{ij} \neq x_{kj}\}$ denote the set of at least $c \cdot d$ indices where x_i and x_k differ, we have

$$\|Z - x_i\|_1 \geq \|Z - x_k\|_1 \quad \text{if and only if} \quad \sum_{j \in J} |Z_j - x_{ij}| - |Z_j - x_{kj}| \geq 0.$$

If x_i is the word being sent and x_i and x_k differ in position j , then $|Z_j - x_{ij}| - |Z_j - x_{kj}| \in \{-1, 1\}$, and is equal to -1 with probability $(1 - \epsilon)$ and 1 with probability ϵ . That is, we have $\|Z - x_i\|_1 \geq \|Z - x_k\|_1$ if and only if

$$\sum_{j \in J} |Z_j - x_{ij}| - |Z_j - x_{kj}| + |J|(1 - 2\epsilon) \geq |J|(1 - 2\epsilon) \geq cd(1 - 2\epsilon),$$

and the expectation $\mathbb{E}_Q[|Z_j - x_{ij}| - |Z_j - x_{kj}| \mid x_i] = -(1 - 2\epsilon)$ when $x_{ij} \neq x_{kj}$. Using the Hoeffding bound, then, we have

$$Q(\|Z - x_i\|_1 \geq \|Z - x_k\|_1 \mid x_i) \leq \exp\left(-\frac{|J|(1 - 2\epsilon)^2}{2}\right) \leq \exp\left(-\frac{cd(1 - 2\epsilon)^2}{2}\right),$$

where we have used that there are at least $|J| \geq cd$ indices differing between x_i and x_k . The probability of making a mistake at all is thus at most $m \exp(-\frac{1}{2}cd(1 - 2\epsilon)^2)$ if our codebook has separation $c \cdot d$.

For low error decoding to occur with extremely high probability, it is thus sufficient to choose a set of code words $\{x_1, \dots, x_m\}$ that is well separated. To that end, we state a simple lemma.

Lemma 4.1.26. *Let X_i , $i = 1, \dots, m$ be drawn independently and uniformly on the d -dimensional hypercube $\mathcal{H}_d := \{0, 1\}^d$. Then for any $t \geq 0$,*

$$\mathbb{P}\left(\exists i, j \text{ s.t. } \|X_i - X_j\|_1 < \frac{d}{2} - dt\right) \leq \binom{m}{2} \exp(-2dt^2) \leq \frac{m^2}{2} \exp(-2dt^2).$$

Proof First, let us consider two independent draws X and X' uniformly on the hypercube. Let $Z = \sum_{j=1}^d \mathbf{1}\{X_j \neq X'_j\} = d_{\text{ham}}(X, X') = \|X - X'\|_1$. Then $\mathbb{E}[Z] = \frac{d}{2}$. Moreover, Z is an i.i.d. sum of Bernoulli $\frac{1}{2}$ random variables, so that by our concentration bounds of Corollary 4.1.10, we have

$$\mathbb{P}\left(\|X - X'\|_1 \leq \frac{d}{2} - t\right) \leq \exp\left(-\frac{2t^2}{d}\right).$$

Using a union bound gives the remainder of the result. \square

Rewriting the lemma slightly, we may take $\delta \in (0, 1)$. Then

$$\mathbb{P}\left(\exists i, j \text{ s.t. } \|X_i - X_j\|_1 < \frac{d}{2} - \sqrt{d \log \frac{1}{\delta} + d \log m}\right) \leq \delta.$$

As a consequence of this lemma, we see two things:

- (i) If $m \leq \exp(d/16)$, or $d \geq 16 \log m$, then taking $\delta \uparrow 1$, there at least exists a codebook $\{x_1, \dots, x_m\}$ of words that are all separated by at least $d/4$, that is, $\|x_i - x_j\|_1 \geq \frac{d}{4}$ for all i, j .
- (ii) By taking $m \leq \exp(d/32)$, or $d \geq 32 \log m$, and $\delta = e^{-d/32}$, then with probability at least $1 - e^{-d/32}$ —exponentially large in d —a randomly drawn codebook has all its entries separated by at least $\|x_i - x_j\|_1 \geq \frac{d}{4}$.

Summarizing, we have the following result: choose a codebook of m codewords x_1, \dots, x_m uniformly at random from the hypercube $\mathcal{H}_d = \{0, 1\}^d$ with

$$d \geq \max\left\{32 \log m, \frac{8 \log \frac{m}{\delta}}{(1 - 2\epsilon)^2}\right\}.$$

Then with probability at least $1 - 1/m$ over the draw of the codebook, the probability we make a mistake in transmission of any given symbol i over the channel Q is at most δ .

4.2 Martingale methods

The next set of tools we consider constitute our first look at argument sbased on *stability*, that is, how quantities that do not change very much when a single observation changes should concentrate. In this case, we would like to understand more general quantities than sample means, developing a few of the basic cools to understand when functions $f(X_1, \dots, X_n)$ of independent random variables X_i concentrate around their expectations. Roughly, we expect that if changing the value of one x_i does not significantly change $f(x_1^n)$ much—it is stable—then it should exhibit good concentration properties.

To develop the tools to do this, we go throughg an approach based on martingales, a deep subject in probability theory. We give a high-level treatment of martingales, taking an approach that does not require measure-theoretic considerations, providing references at the end of the chapter. We begin by providing a definition.

Definition 4.3. Let M_1, M_2, \dots be an \mathbb{R} -valued sequence of random variables. They are a martingale if there exist another sequence of random variables $\{Z_1, Z_2, \dots\} \subset \mathcal{Z}$ and sequence of functions $f_n : \mathcal{Z}^n \rightarrow \mathbb{R}$ such that

$$\mathbb{E}[M_n \mid Z_1^{n-1}] = M_{n-1} \quad \text{and} \quad M_n = f_n(Z_1^n)$$

for all $n \in \mathbb{N}$. We say that the sequence M_n is adapted to $\{Z_n\}$.

In general, the sequence Z_1, Z_2, \dots is a sequence of increasing σ -fields $\mathcal{F}_1, \mathcal{F}_2, \dots$, and M_n is \mathcal{F}_n -measurable, but Definition 4.3 is sufficienet for our purposes. We also will find it convenient to study *differences* of martingales, so that we make the following

Definition 4.4. Let D_1, D_2, \dots be a sequence of random variables. They form a martingale difference sequence if $M_n := \sum_{i=1}^n D_i$ is a martingale.

Equivalently, there is a sequence of random variables Z_n and functions $g_n : \mathcal{Z}^n \rightarrow \mathbb{R}$ such that

$$\mathbb{E}[D_n \mid Z_1^{n-1}] = 0 \quad \text{and} \quad D_n = g_n(Z_1^n)$$

for all $n \in \mathbb{N}$.

There are numerous examples of martingale sequences. The classical one is the symmetric random walk.

Example 4.2.1: Let $D_n \in \{\pm 1\}$ be uniform and independent. Then D_n form a martingale difference sequence adapted to themselves (that is, we may take $Z_n = D_n$), and $M_n = \sum_{i=1}^n D_i$ is a martingale. \diamond

A more sophisticated example, to which we will frequently return and that suggests the potential usefulness of martingale constructions, is the *Doob martingale* associated with a function f .

Example 4.2.2 (Doob martingales): Let $f : \mathcal{X}^n \rightarrow \mathbb{R}$ be an otherwise arbitrary function, and let X_1, \dots, X_n be arbitrary random variables. The Doob martingale is defined by the difference sequence

$$D_i := \mathbb{E}[f(X_1^n) \mid X_1^i] - \mathbb{E}[f(X_1^n) \mid X_1^{i-1}].$$

By inspection, the D_i are functions of X_1^i , and we have

$$\begin{aligned} \mathbb{E}[D_i \mid X_1^{i-1}] &= \mathbb{E}[\mathbb{E}[f(X_1^n) \mid X_1^i] \mid X_1^{i-1}] - \mathbb{E}[f(X_1^n) \mid X_1^{i-1}] \\ &= \mathbb{E}[f(X_1^n) \mid X_1^{i-1}] - \mathbb{E}[f(X_1^n) \mid X_1^{i-1}] = 0 \end{aligned}$$

by the tower property of expectations. Thus, the D_i satisfy Definition 4.4 of a martingale difference sequence, and moreover, we have

$$\sum_{i=1}^n D_i = f(X_1^n) - \mathbb{E}[f(X_1^n)],$$

and so the Doob martingale captures exactly the difference between f and its expectation. \diamond

4.2.1 Sub-Gaussian martingales and Azuma-Hoeffding inequalities

With these motivating ideas introduced, we turn to definitions, providing generalizations of our concentration inequalities for sub-Gaussian sums to sub-Gaussian martingales, which we define.

Definition 4.5. Let $\{D_n\}$ be a martingale difference sequence adapted to $\{Z_n\}$. Then D_n is a σ_n^2 -sub-Gaussian martingale difference if

$$\mathbb{E}[\exp(\lambda D_n) \mid Z_1^{n-1}] \leq \exp\left(\frac{\lambda^2 \sigma_n^2}{2}\right)$$

for all n and $\lambda \in \mathbb{R}$.

Immediately from the definition, we have the Azuma-Hoeffding inequalities, which generalize the earlier tensorization identities for sub-Gaussian random variables.

Theorem 4.2.3 (Azuma-Hoeffding). Let $\{D_n\}$ be a σ_n^2 -sub-Gaussian martingale difference sequence. Then $M_n = \sum_{i=1}^n D_i$ is $\sum_{i=1}^n \sigma_i^2$ -sub-Gaussian, and moreover,

$$\max\{\mathbb{P}(M_n \geq t), \mathbb{P}(M_n \leq -t)\} \leq \exp\left(-\frac{nt^2}{2 \sum_{i=1}^n \sigma_i^2}\right) \text{ for all } t \geq 0.$$

Proof The proof is essentially immediate: letting Z_n be the sequence to which the D_n are adapted, we write

$$\begin{aligned} \mathbb{E}[\exp(\lambda M_n)] &= \mathbb{E}\left[\prod_{i=1}^n e^{\lambda D_i}\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[\prod_{i=1}^n e^{\lambda D_i} \mid Z_1^{n-1}\right]\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[\prod_{i=1}^{n-1} e^{\lambda D_i} \mid Z_1^{n-1}\right] \mathbb{E}[e^{\lambda D_n} \mid Z_1^{n-1}]\right] \end{aligned}$$

because D_1, \dots, D_{n-1} are functions of Z_1^{n-1} . Then we use Definition 4.5, which implies that $\mathbb{E}[e^{\lambda D_n} \mid Z_1^{n-1}] \leq e^{\lambda^2 \sigma_n^2 / 2}$, and we obtain

$$\mathbb{E}[\exp(\lambda M_n)] \leq \mathbb{E}\left[\prod_{i=1}^{n-1} e^{\lambda D_i}\right] \exp\left(\frac{\lambda^2 \sigma_n^2}{2}\right).$$

Repeating the same argument for $n-1, n-2, \dots, 1$ gives that

$$\log \mathbb{E}[\exp(\lambda M_n)] \leq \frac{\lambda^2}{2} \sum_{i=1}^n \sigma_i^2$$

as desired.

The second claims are simply applications of Chernoff bounds via Proposition 4.1.8 and that $\mathbb{E}[M_n] = 0$. \square

As an immediate corollary, we recover Proposition 4.1.9, as sums of independent random variables form martingales via $M_n = \sum_{i=1}^n (X_i - \mathbb{E}[X_i])$. A second corollary gives what is typically termed the Azuma inequality:

Corollary 4.2.4. *Let D_i be a bounded difference martingale difference sequence, meaning that $|D_i| \leq c$. Then $M_n = \sum_{i=1}^n D_i$ satisfies*

$$\mathbb{P}(n^{-1/2}M_n \geq t) \vee \mathbb{P}(n^{-1/2}M_n \leq -t) \leq \exp\left(-\frac{t^2}{2c^2}\right) \quad \text{for } t \geq 0.$$

Thus, bounded random walks are (with high probability) within $\pm\sqrt{n}$ of their expectations after n steps.

There exist extensions of these inequalities to the cases where we control the variance of the martingales; see Freedman [96].

4.2.2 Examples and bounded differences

We now develop several example applications of the Azuma-Hoeffding inequalities (Theorem 4.2.3), applying them most specifically to functions satisfying certain stability conditions.

We first define the collections of functions we consider.

Definition 4.6 (Bounded differences). *Let $f : \mathcal{X}^n \rightarrow \mathbb{R}$ for some space \mathcal{X} . Then f satisfies bounded differences with constants c_i if for each $i \in \{1, \dots, n\}$, all $x_1^n \in \mathcal{X}^n$, and $x'_i \in \mathcal{X}$ we have*

$$|f(x_1^{i-1}, x_i, x_{i+1}^n) - f(x_1^{i-1}, x'_i, x_{i+1}^n)| \leq c_i.$$

The classical inequality relating bounded differences and concentration is McDiarmid's inequality, or the bounded differences inequality.

Proposition 4.2.5 (Bounded differences inequality). *Let $f : \mathcal{X}^n \rightarrow \mathbb{R}$ satisfy bounded differences with constants c_i , and let X_i be independent random variables. $f(X_1^n) - \mathbb{E}[f(X_1^n)]$ is $\frac{1}{4} \sum_{i=1}^n c_i^2$ -sub-Gaussian, and*

$$\mathbb{P}(f(X_1^n) - \mathbb{E}[f(X_1^n)] \geq t) \vee \mathbb{P}(f(X_1^n) - \mathbb{E}[f(X_1^n)] \leq -t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right).$$

Proof The basic idea is to show that the Doob martingale (Example 4.2.2) associated with f is $c_i^2/4$ -sub-Gaussian, and then to simply apply the Azuma-Hoeffding inequality. To that end, define $D_i = \mathbb{E}[f(X_1^n) | X_i] - \mathbb{E}[f(X_1^n) | X_1^{i-1}]$ as before, and note that $\sum_{i=1}^n D_i = f(X_1^n) - \mathbb{E}[f(X_1^n)]$. The random variables

$$\begin{aligned} L_i &:= \inf_x \mathbb{E}[f(X_1^n) | X_1^{i-1}, X_i = x] - \mathbb{E}[f(X_1^n) | X_1^{i-1}] \\ U_i &:= \sup_x \mathbb{E}[f(X_1^n) | X_1^{i-1}, X_i = x] - \mathbb{E}[f(X_1^n) | X_1^{i-1}] \end{aligned}$$

evidently satisfy $L_i \leq D_i \leq U_i$, and moreover, we have

$$\begin{aligned} U_i - L_i &\leq \sup_{x_1^{i-1}} \sup_{x, x'} \{ \mathbb{E}[f(X_1^n) \mid X_1^{i-1} = x_1^{i-1}, X_i = x] - \mathbb{E}[f(X_1^n) \mid X_1^{i-1} = x_1^{i-1}, X_i = x'] \} \\ &= \sup_{x_1^{i-1}} \sup_{x, x'} \int (f(x_1^{i-1}, x, x_{i+1}^n) - f(x_1^{i-1}, x', x_{i+1}^n)) dP(x_{i+1}^n) \leq c_i, \end{aligned}$$

where we have used the independence of the X_i and Definition 4.6 of bounded differences. Consequently, we have by Hoeffding's Lemma (Example 4.1.6) that $\mathbb{E}[e^{\lambda D_i} \mid X_1^{i-1}] \leq \exp(\lambda^2 c_i^2 / 8)$, that is, the Doob martingale is $c_i^2 / 4$ -sub-Gaussian.

The remainder of the proof is simply Theorem 4.2.3. \square

A number of quantities satisfy the conditions of Proposition 4.2.5, and we give two examples here; we will revisit them more later.

Example 4.2.6 (Bounded random vectors): Let \mathbb{B} be a Banach space—a complete normed vector space—with norm $\|\cdot\|$. Let X_i be independent bounded random vectors in \mathbb{B} satisfying $\mathbb{E}[X_i] = 0$ and $\|X_i\| \leq c$. We claim that the quantity

$$f(X_1^n) := \left\| \frac{1}{n} \sum_{i=1}^n X_i \right\|$$

satisfies bounded differences. Indeed, we have by the triangle inequality that

$$|f(x_1^{i-1}, x, x_{i+1}^n) - f(x_1^{i-1}, x', x_{i+1}^n)| \leq \frac{1}{n} \|x - x'\| \leq \frac{2c}{n}.$$

Consequently, if X_i are independent, we have

$$\mathbb{P} \left(\left| \left\| \frac{1}{n} \sum_{i=1}^n X_i \right\| - \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n X_i \right\| \right] \right| \geq t \right) \leq 2 \exp \left(-\frac{nt^2}{2c^2} \right) \quad (4.2.1)$$

for all $t \geq 0$. That is, the norm of (bounded) random vectors in an essentially arbitrary vector space concentrates extremely quickly about its expectation.

The challenge becomes to control the *expectation* term in the concentration bound (4.2.1), which can be a bit challenging. In certain cases—for example, when we have a Euclidean structure on the vectors X_i —it can be easier. Indeed, let us specialize to the case that $X_i \in \mathcal{H}$, a (real) Hilbert space, so that there is an inner product $\langle \cdot, \cdot \rangle$ and the norm satisfies $\|x\|^2 = \langle x, x \rangle$ for $x \in \mathcal{H}$. Then Cauchy-Schwarz implies that

$$\mathbb{E} \left[\left\| \sum_{i=1}^n X_i \right\|^2 \right] \leq \mathbb{E} \left[\left\| \sum_{i=1}^n X_i \right\|^2 \right] = \sum_{i,j} \mathbb{E}[\langle X_i, X_j \rangle] = \sum_{i=1}^n \mathbb{E}[\|X_i\|^2].$$

That is assuming the X_i are independent and $\mathbb{E}[\|X_i\|^2] \leq \sigma^2$, inequality (4.2.1) becomes

$$\mathbb{P} \left(\|\bar{X}_n\| \geq \frac{\sigma}{\sqrt{n}} + t \right) + \mathbb{P} \left(\|\bar{X}_n\| \leq -\frac{\sigma}{\sqrt{n}} - t \right) \leq 2 \exp \left(-\frac{nt^2}{2c^2} \right)$$

where $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. \diamond

We can specialize Example 4.2.6 to a situation that is very important for treatments of concentration, sums of random vectors, and generalization bounds in machine learning.

Example 4.2.7 (Rademacher complexities): This example is actually a special case of Example 4.2.6, but its frequent uses justify a more specialized treatment and consideration. Let \mathcal{X} be some space, and let \mathcal{F} be some collection of functions $f : \mathcal{X} \rightarrow \mathbb{R}$. Let $\varepsilon_i \in \{-1, 1\}$ be a collection of independent random sign vectors. Then the *empirical Rademacher complexity* of \mathcal{F} is

$$R_n(\mathcal{F} \mid x_1^n) := \frac{1}{n} \mathbb{E} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^n \varepsilon_i f(x_i) \right],$$

where the expectation is over only the random signs ε_i . (In some cases, depending on context and convenience, one takes the absolute value $|\sum_i \varepsilon_i f(x_i)|$.) The *Rademacher complexity* of \mathcal{F} is

$$R_n(\mathcal{F}) := \mathbb{E}[R_n(\mathcal{F} \mid X_1^n)],$$

the expectation of the empirical Rademacher complexities.

If $f : \mathcal{X} \rightarrow [b_0, b_1]$ for all $f \in \mathcal{F}$, then the Rademacher complexity satisfies bounded differences, because for any two sequences x_1^n and z_1^n differing in only element j , we have

$$n|R_n(\mathcal{F} \mid x_1^n) - R_n(\mathcal{F} \mid z_1^n)| \leq \mathbb{E} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^n \varepsilon_i (f(x_i) - f(z_i)) \right] = \mathbb{E}[\sup_{f \in \mathcal{F}} \varepsilon_i (f(x_j) - f(z_j))] \leq b_1 - b_0.$$

Consequently, the empirical Rademacher complexity satisfies $R_n(\mathcal{F} \mid X_1^n) - R_n(\mathcal{F})$ is $\frac{(b_1 - b_0)^2}{4n}$ -sub-Gaussian by Theorem 4.2.3. \diamond

These examples warrant more discussion, and it is possible to argue that many variants of these random variables are well-concentrated. For example, instead of functions we may simply consider an arbitrary set $\mathcal{A} \subset \mathbb{R}^n$ and define the random variable

$$Z(\mathcal{A}) := \sup_{a \in \mathcal{A}} \langle a, \varepsilon \rangle = \sup_{a \in \mathcal{A}} \sum_{i=1}^n a_i \varepsilon_i.$$

As a function of the random signs ε_i , we may write $Z(\mathcal{A}) = f(\varepsilon)$, and this is then a function satisfying $|f(\varepsilon) - f(\varepsilon')| \leq \sup_{a \in \mathcal{A}} |\langle a, \varepsilon - \varepsilon' \rangle|$, so that if ε and ε' differ in index i , we have $|f(\varepsilon) - f(\varepsilon')| \leq 2 \sup_{a \in \mathcal{A}} |a_i|$. That is, $Z(\mathcal{A}) - \mathbb{E}[Z(\mathcal{A})]$ is $\sum_{i=1}^n \sup_{a \in \mathcal{A}} |a_i|^2$ -sub-Gaussian.

Example 4.2.8 (Rademacher complexity as a random vector): This view of Rademacher complexity shows how we may think of Rademacher complexities as norms on certain spaces. Indeed, if we consider a vector space \mathcal{L} of linear functions on \mathcal{F} , then we can define the \mathcal{F} -seminorm on \mathcal{L} by $\|L\|_{\mathcal{F}} := \sup_{f \in \mathcal{F}} |L(f)|$. In this case, we may consider the symmetrized empirical distributions

$$P_n^0 := \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{1}_{X_i} \quad f \mapsto P_n^0 f := \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i)$$

as elements of this vector space \mathcal{L} . (Here we have used $\mathbf{1}_{X_i}$ to denote the point mass at X_i .) Then the Rademacher complexity is nothing more than the expected norm of P_n^0 , a random vector, as in Example 4.2.6. This view is somewhat sophisticated, but it shows that any general results we may prove about random vectors, as in Example 4.2.6, will carry over immediately to versions of the Rademacher complexity. \diamond

4.3 Matrix concentration

In this section, we will develop analogues of the concentration inequalities for sums in Section 4.1, including matrix Hoeffding and Bernstein inequalities. Our main goal will be to bound maximal eigenvalues (or operator norms) of symmetric and Hermitian matrices, that is, for sums $S_n = \sum_{i=1}^n X_i$ of independent matrices, of deviation probabilities

$$\mathbb{P}(\lambda_{\max}(S_n) \geq t) \text{ or } \mathbb{P}(\lambda_{\min}(S_n) \leq -t),$$

where λ_{\max} and λ_{\min} denote maximal and minimal eigenvalues, respectively. Our approach will be to generalize the approach using moment generating functions, though this becomes non-trivial because there is no immediately obvious analogue of the tensorization identities we have for scalars. While in the scalar case, for a sum $S_n = \sum_{i=1}^n X_i$ of independent random variables, we have

$$e^{\lambda S_n} = \prod_{i=1}^n e^{\lambda X_i},$$

such an identity fails for matrices, because their exponentials (typically) fail to commute.

To develop the basic matrix concentration equalities we provide, we require a brief review of matrix calculus and operator functions. We shall typically work with Hermitian matrices $A \in \mathbb{C}^{d \times d}$, meaning that $A = A^*$, where A^* denotes the Hermitian transpose of A , whose entries are $(A^*)_{ij} = \overline{A_{ji}}$, the conjugate of A_{ji} . We work in this generality for two reasons: first, because such matrices admit the spectral decompositions we require to develop the operators we use, and second, because we often will encounter random matrices with symmetric distributions, meaning that $X \stackrel{\text{dist}}{=} -X$, which can lead to confusion.

With this, we give a brief review of some properties of Hermitian matrices and some associated matrix operators. Let $\mathcal{H}_d := \{A \in \mathbb{C}^{d \times d} \mid A^* = A\}$ be the Hermitian matrices. The spectral theorem states gives that any $A \in \mathcal{H}_d$ admits the spectral decomposition $A = U \Lambda U^*$, where Λ is the diagonal matrix of the (necessarily) real eigenvalues of A and $U \in \mathbb{C}^{n \times n}$ is unitary, so that $U^*U = UU^* = I$. For a function $f : \mathbb{R} \rightarrow \mathbb{R}$, we can then define its *operator extension* to \mathcal{H}_d by

$$f(A) := U \text{diag}(f(\lambda_1(A)), \dots, f(\lambda_n(A))) U^*,$$

where A has spectral decomposition $A = U \Lambda U^*$ and $\lambda_i(A)$ denotes the i th eigenvalue of A . Because we wish to mimic the approach based on moment generating functions that yields our original sub-Gaussian and sub-Exponential concentration inequalities in Chapter 4, the most important function for us will be the exponential, which evidently satisfies

$$\exp(A) = \sum_{k=0}^{\infty} \frac{1}{k!} A^k,$$

where we recall the convention that $A^0 = I$ whenever A is Hermitian.

A Hermitian matrix A is *positive definite*, denoted $A \succ 0$, if $x^* A x > 0$ for all $x \neq 0$, and is positive semidefinite (PSD), which we denote by $A \succeq 0$, if $x^* A x \geq 0$ for all vectors x . Positive definiteness is then equivalent to the condition that $\lambda_i(A) > 0$ for all eigenvalues of A , while semidefiniteness that $\lambda_i(A) \geq 0$. We also use the standard semidefinite ordering, so that $A \succeq B$ means that $A - B \succeq 0$. For $A \in \mathcal{H}_d$, we evidently have $\exp(A) \succ 0$. The familiar trace $\text{tr}(A) = \sum_{j=1}^d A_{jj}$ of a square matrix allows us to define inner products, where for general complex matrices

$A, B \in \mathbb{C}^{m \times n}$ we define $\langle A, B \rangle = \text{tr}(A^*B)$, while the space of Hermitian matrices admits the *real* inner product $\langle A, B \rangle := \text{tr}(A^*B)$. (See Exercise 4.14.) The spectral theorem also shows the standard identity that $\text{tr}(A) = \sum_{j=1}^d \lambda_j(A)$ for $A \in \mathcal{H}_d$.

To analogize our approach with real-valued random variables, we begin with the Chernoff bound, Proposition 4.1.3. Here, we have the following observation:

Proposition 4.3.1. *For any random Hermitian matrix X ,*

$$\mathbb{P}(\lambda_{\max}(X) \geq t) \leq \text{tr}(\mathbb{E}[e^{\lambda X}])e^{-\lambda t}$$

for all $\lambda \geq 0$ and $t \geq 0$.

Proof First, apply the standard Chernoff bound to the random variable $\lambda_{\max}(X)$, which gives that for any $\lambda > 0$ that

$$\mathbb{P}(\lambda_{\max}(X) \geq t) \leq \mathbb{E}[e^{\lambda \lambda_{\max}(X)}]e^{-\lambda t}.$$

Then observe that by definition of the matrix exponential, we have $e^{\lambda \lambda_{\max}(X)} \leq \text{tr}(e^{\lambda X})$, because the eigenvalues of $e^{\lambda X}$ are all positive. \square

We would like now to provide some type of general tensorization identity for matrices, in analogy with Propositions 4.1.9 or 4.1.17. Unfortunately, this breaks down: for Hermitian A, B , we have

$$e^{A+B} = e^A e^B$$

if and only if A and B commute [153], so that they are simultaneously diagonalizable. Nonetheless, we have the following inequality, which will be the key to extending the standard one-dimensional approach to concentration:

Proposition 4.3.2 (The Golden-Thompson inequality). *Let A, B be Hermitian matrices. Then*

$$\text{tr}(e^{A+B}) \leq \text{tr}(e^A e^B).$$

While the proof is essentially elementary, it is not central to our development, so we defer it to Section 4.4.4. We remark in passing that there is a converse [153, Section 3]: $\text{tr}(e^{A+B}) = \text{tr}(e^A e^B)$ if and only if $AB = BA$, that is, A and B are simultaneously diagonalizable. With Proposition 4.3.2 in hand, however, we can develop matrix analogues of the Hoeffding and Bernstein-type concentration bounds in Chapter 4.

We begin with Azuma-Hoeffding-type bounds, which analogize Theorem 4.2.3. The key to allow an iterative “peeling” off of individual terms in a sum of random matrices is the following result:

Lemma 4.3.3 (A matrix symmetrization inequality). *Let H be an arbitrary (fixed) Hermitian matrix and X be a mean-zero Hermitian matrix. Then*

$$\text{tr}(\mathbb{E}[e^{H+X}]) \leq \text{tr}(\mathbb{E}[e^{H+2\varepsilon X}]).$$

Proof Let X' be an independent copy of X . Then because the trace exponential $\text{tr}(e^X)$ is convex on the Hermitian matrices (see Exercise 4.16), we have

$$\text{tr}(\mathbb{E}[e^{H+X}]) = \text{tr}(\mathbb{E}[e^{H+(X-\mathbb{E}[X'])}]) \leq \text{tr}(\mathbb{E}[e^{H+X-X'}])$$

by Jensen's inequality. Introducing the random sign $\varepsilon \in \{\pm 1\}$, we have by symmetry that $X - X' \stackrel{\text{dist}}{=} \varepsilon(X - X')$, and so

$$\begin{aligned} \text{tr}(\mathbb{E}[e^{H+X}]) &\leq \mathbb{E}[\text{tr}(e^{H+\varepsilon X - \varepsilon X'})] = \mathbb{E}[\text{tr}(e^{H/2+\varepsilon X + H/2-\varepsilon X'})] \\ &\leq \mathbb{E}[\text{tr}(e^{H/2+\varepsilon X} e^{H/2-\varepsilon X'})], \end{aligned}$$

where the second inequality follows from Proposition 4.3.2. Now we use that for Hermitian matrices A, B , we have $\text{tr}(AB) \leq \|A\|_2 \|B\|_2 = \text{tr}(A^2)^{1/2} \text{tr}(B^2)^{1/2}$, so that

$$\mathbb{E}[\text{tr}(e^{H/2+\varepsilon X} e^{H/2-\varepsilon X'})] \leq \mathbb{E}[\text{tr}(e^{H+2\varepsilon X})^{1/2} \text{tr}(e^{H-2\varepsilon X'})^{1/2}] \leq \mathbb{E}[\text{tr}(e^{H+2\varepsilon X})]^{1/2} \mathbb{E}[\text{tr}(e^{H-2\varepsilon X'})]^{1/2}$$

by Cauchy-Schwarz. Because $\varepsilon X \stackrel{\text{dist}}{=} -\varepsilon X'$, the lemma follows. \square

This allows us to perform the type of “peeling-off” argument, addressing one term in the sum at a type, that gives tight enough moment generating function bounds.

Theorem 4.3.4. *Let $X_1, \dots, X_n \in \mathcal{H}_d$ be independent, mean-zero, and satisfy $\|X_i\|_{\text{op}} \leq b_i$. Define $S_n = \sum_{i=1}^n X_i$. Then for all $\lambda \in \mathbb{R}$,*

$$\mathbb{E}[\text{tr}(e^{\lambda S_n})] \leq d \exp \left(2\lambda^2 \sum_{i=1}^n b_i^2 \right).$$

Proof By iterated expectation and Lemma 4.3.3, we have

$$\text{tr}(\mathbb{E}[e^{\lambda S_n}]) \leq \text{tr}(\mathbb{E}[e^{\lambda S_{n-1} + 2\lambda \varepsilon X_n}]) \leq \text{tr}(\mathbb{E}[e^{\lambda S_{n-1}}]) \left\| \mathbb{E}[e^{2\lambda \varepsilon X_n}] \right\|_{\text{op}},$$

where $\varepsilon \in \{\pm 1\}$ is an independent random sign and we have used independence. Now, we use the following calculation: if X is Hermitian and ε a random sign, then

$$\mathbb{E}[e^{\varepsilon X}] \preceq \mathbb{E}[e^{X^2/2}]. \quad (4.3.1)$$

Temporarily deferring the argument for inequality (4.3.1), note that it immediately implies $\mathbb{E}[e^{2\lambda \varepsilon X}] \preceq \mathbb{E}[e^{2\lambda^2 X^2}]$. The convexity of the operator norm and that $\|X_n\|_{\text{op}} \leq b_n$ then imply

$$\left\| \mathbb{E}[e^{2\lambda \varepsilon X_n}] \right\|_{\text{op}} \leq \left\| \mathbb{E}[e^{2\lambda^2 X_n^2}] \right\|_{\text{op}} \leq \mathbb{E} \left[e^{2\lambda^2 \|X_n\|_{\text{op}}^2} \right] \leq e^{2\lambda^2 b_n^2}$$

Repeating the argument by iteratively peeling off the last term X_{n-i} for S_{n-2} through S_1 then yields

$$\mathbb{E}[\text{tr}(\exp(\lambda S_n))] \leq \text{tr}(I) \prod_{i=1}^n \exp(2\lambda^2 b_i^2),$$

which gives the theorem.

To see inequality (4.3.1), note that for any positive semidefinite A , we have $A \preceq tA$ for $t \geq 1$. Then because $X^{2k} \succeq 0$ for all $k \in \mathbb{N}$ and $(2k)! \geq 2^k k!$, we have

$$\mathbb{E}[e^{\varepsilon X}] = I + \sum_{k=1}^{\infty} \frac{\mathbb{E}[X^{2k}]}{(2k)!} \preceq I + \sum_{k=1}^{\infty} \frac{\mathbb{E}[(X^2)^k]}{2^k k!} = \mathbb{E}[e^{X^2/2}],$$

where we used symmetry to eliminate terms with odd powers. \square

Theorem 4.3.4 immediately implies the following corollary, whose argument parallels those in Chapter 4 (e.g., Corollary 4.1.10).

Corollary 4.3.5. *Let $X_i \in \mathcal{H}_d$ be independent mean-zero Hermitian matrices. Then $S_n := \sum_{i=1}^n X_i$ satisfies*

$$\mathbb{P}\left(\|S_n\|_{\text{op}} \geq t\right) \leq 2d \exp\left(-\frac{t^2}{8 \sum_{i=1}^n b_i^2}\right).$$

If we have more direct bounds on $\mathbb{E}[e^{\lambda X_i}]$, then we can also employ those via a similar “peeling off” the last term argument. By carefully controlling matrix moment generating functions in a way similar to that we did in Example 4.1.14 to obtain sub-exponential behavior for bounded random variables, we can give a matrix Bernstein-type inequality.

Theorem 4.3.6. *Let X_i be independent Hermitian matrices with $\|X_i\|_{\text{op}} \leq b$ and $\|\mathbb{E}[X_i^2]\|_{\text{op}} \leq \sigma_i^2$. Then $S_n = \sum_{i=1}^n X_i$ satisfies*

$$\mathbb{P}\left(\|S_n\|_{\text{op}} \geq t\right) \leq 2d \exp\left(-\min\left\{\frac{t^2}{4 \sum_{i=1}^n \sigma_i^2}, \frac{3t}{4b}\right\}\right).$$

The proof of the theorem is similar to that of Theorem 4.3.4, so we leave it as an extended exercise (Exercise 4.17).

We unpack the theorem a bit to give some intuition. Given a variance bound σ^2 such that $\mathbb{E}[X_i^2] \preceq \sigma^2 I$, the theorem states that

$$\mathbb{P}\left(\|n^{-1}S_n\|_{\text{op}} \geq t\right) \leq 2d \exp\left(-\min\left\{\frac{nt^2}{4\sigma^2}, \frac{3nt}{4b}\right\}\right).$$

Letting $\delta \in (0, 1)$ be arbitrary and setting $t = \max\left\{\frac{2\sigma}{\sqrt{n}}\sqrt{\log \frac{2d}{\delta}}, \frac{4b}{3n} \log \frac{d}{\delta}\right\}$, we have

$$\left\|\frac{1}{n} \sum_{i=1}^n X_i\right\|_{\text{op}} \leq \max\left\{\frac{2\sigma}{\sqrt{n}}\sqrt{\log \frac{2d}{\delta}}, \frac{4b}{3n} \log \frac{d}{\delta}\right\}$$

with probability at least $1 - \delta$. So we see the familiar sub-Gaussian and sub-exponential scaling of the random sum.

4.4 Technical proofs

4.4.1 Proof of Theorem 4.1.11

(1) **implies** (2) Let $K_1 = 1$. Using the change of variables identity that for a nonnegative random variable Z and any $k \geq 1$ we have $\mathbb{E}[Z^k] = k \int_0^\infty t^{k-1} \mathbb{P}(Z \geq t) dt$, we find

$$\mathbb{E}[|X|^k] = k \int_0^\infty t^{k-1} \mathbb{P}(|X| \geq t) dt \leq 2k \int_0^\infty t^{k-1} \exp\left(-\frac{t^2}{\sigma^2}\right) dt = k\sigma^k \int_0^\infty u^{k/2-1} e^{-u} du,$$

where for the last inequality we made the substitution $u = t^2/\sigma^2$. Noting that this final integral is $\Gamma(k/2)$, we have $\mathbb{E}[|X|^k] \leq k\sigma^k \Gamma(k/2)$. Because $\Gamma(s) \leq s^s$ for $s \geq 1$, we obtain

$$\mathbb{E}[|X|^k]^{1/k} \leq k^{1/k} \sigma \sqrt{k/2} \leq e^{1/e} \sigma \sqrt{k}.$$

Thus (2) holds with $K_2 = e^{1/e}$.

(2) **implies** (3) Let $\sigma = \sup_{k \geq 1} k^{-\frac{1}{2}} \mathbb{E}[|X|^k]^{1/k}$, so that $K_2 = 1$ and $\mathbb{E}[|X|^k] \leq k^{\frac{k}{2}} \sigma$ for all k . For $K_3 \in \mathbb{R}_+$, we thus have

$$\mathbb{E}[\exp(X^2/(K_3\sigma^2))] = \sum_{k=0}^{\infty} \frac{\mathbb{E}[X^{2k}]}{k! K_3^{2k} \sigma^{2k}} \leq \sum_{k=0}^{\infty} \frac{\sigma^{2k} (2k)^k}{k! K_3^{2k} \sigma^{2k}} \stackrel{(i)}{\leq} \sum_{k=0}^{\infty} \left(\frac{2e}{K_3^2} \right)^k$$

where inequality (i) follows because $k! \geq (k/e)^k$, or $1/k! \leq (e/k)^k$. Noting that $\sum_{k=0}^{\infty} \alpha^k = \frac{1}{1-\alpha}$, we obtain (3) by taking $K_3 = e\sqrt{2/(e-1)} \approx 2.933$.

(3) **implies** (4) Let us take $K_3 = 1$ and recall the assumption of (4) that $\mathbb{E}[X] = 0$. We claim that (4) holds with $K_4 = \frac{3}{4}$. We prove this result for both small and large λ . First, note the (non-standard but true!) inequality that $e^x \leq x + e^{\frac{9x^2}{16}}$ for all x . Then we have

$$\mathbb{E}[\exp(\lambda X)] \leq \underbrace{\mathbb{E}[\lambda X]}_{=0} + \mathbb{E} \left[\exp \left(\frac{9\lambda^2 X^2}{16} \right) \right]$$

Now note that for $|\lambda| \leq \frac{4}{3\sigma}$, we have $9\lambda^2\sigma^2/16 \leq 1$, and so by Jensen's inequality,

$$\mathbb{E} \left[\exp \left(\frac{9\lambda^2 X^2}{16} \right) \right] = \mathbb{E} \left[\exp(X^2/\sigma^2)^{\frac{9\lambda^2\sigma^2}{16}} \right] \leq e^{\frac{9\lambda^2\sigma^2}{16}}.$$

For large λ , we use the simpler Fenchel-Young inequality, that is, that $\lambda x \leq \frac{\lambda^2}{2c} + \frac{cx^2}{2}$, valid for all $c \geq 0$. Then we have for any $0 \leq c \leq 2$ that

$$\mathbb{E}[\exp(\lambda X)] \leq e^{\frac{\lambda^2\sigma^2}{2c}} \mathbb{E} \left[\exp \left(\frac{cX^2}{2\sigma^2} \right) \right] \leq e^{\frac{\lambda^2\sigma^2}{2c}} e^{\frac{c}{2}},$$

where the final inequality follows from Jensen's inequality. If $|\lambda| \geq \frac{4}{3\sigma}$, then $\frac{1}{2} \leq \frac{9}{32}\lambda^2\sigma^2$, and we have

$$\mathbb{E}[\exp(\lambda X)] \leq \inf_{c \in [0, 2]} e^{\left[\frac{1}{2c} + \frac{9c}{32}\right]\lambda^2\sigma^2} = \exp \left(\frac{3\lambda^2\sigma^2}{4} \right).$$

(3) **implies** (1) Assume (3) holds with $K_3 = 1$. Then for $t \geq 0$ we have

$$\mathbb{P}(|X| \geq t) = \mathbb{P}(X^2/\sigma^2 \geq t^2/\sigma^2) \leq \mathbb{E}[\exp(\lambda X^2/\sigma^2)] \exp \left(-\frac{\lambda t^2}{\sigma^2} \right)$$

for all $\lambda \geq 0$. For $\lambda \leq 1$, Jensen's inequality implies $\mathbb{E}[\exp(\lambda X^2/\sigma^2)] \leq \mathbb{E}[\exp(X^2/\sigma^2)]^\lambda \leq e^\lambda$ by assumption (3). Set $\lambda = \log 2 \approx .693$.

(4) **implies** (1) This is the content of Proposition 4.1.8, with $K_4 = \frac{1}{2}$ and $K_1 = 2$.

4.4.2 Proof of Theorem 4.1.15

(1) **implies** (2) As in the proof of Theorem 4.1.11, we use that for a nonnegative random variable Z we have $\mathbb{E}[Z^k] = k \int_0^\infty t^{k-1} \mathbb{P}(Z \geq t) dt$. Let $K_1 = 1$. Then

$$\mathbb{E}[|X|^k] = k \int_0^\infty t^{k-1} \mathbb{P}(|X| \geq t) dt \leq 2k \int_0^\infty t^{k-1} \exp(-t/\sigma) dt = 2k\sigma^k \int_0^\infty u^{k-1} \exp(-u) du,$$

where we used the substitution $u = t/\sigma$. Thus we have $\mathbb{E}[|X|^k] \leq 2\Gamma(k+1)\sigma^k$, and using $\Gamma(k+1) \leq k^k$ yields $\mathbb{E}[|X|^k]^{1/k} \leq 2^{1/k} k\sigma$, so that (2) holds with $K_2 \leq 2$.

(2) **implies** (3) Let $K_2 = 1$, and note that

$$\mathbb{E}[\exp(X/(K_3\sigma))] = \sum_{k=0}^{\infty} \frac{\mathbb{E}[X^k]}{K_3^k \sigma^k k!} \leq \sum_{k=0}^{\infty} \frac{k^k}{k!} \cdot \frac{1}{K_3^k} \stackrel{(i)}{\leq} \sum_{k=0}^{\infty} \left(\frac{e}{K_3}\right)^k,$$

where inequality (i) used that $k! \geq (k/e)^k$. Taking $K_3 = e^2/(e-1) < 5$ gives the result.

(3) **implies** (1) If $\mathbb{E}[\exp(X/\sigma)] \leq e$, then for $t \geq 0$

$$\mathbb{P}(X \geq t) \leq \mathbb{E}[\exp(X/\sigma)] e^{-t/\sigma} \leq e^{1-t/\sigma}.$$

With the same result for the negative tail, we have

$$\mathbb{P}(|X| \geq t) \leq 2e^{1-t/\sigma} \wedge 1 \leq 2e^{-\frac{2t}{5\sigma}},$$

so that (1) holds with $K_1 = \frac{5}{2}$.

(2) **if and only if** (4) Assume that (2) holds with $K_2 = 1$, and let $\sigma = \sup_{k \geq 1} \frac{1}{k} \mathbb{E}[|X|^k]^{1/k}$. Then because $\mathbb{E}[X] = 0$,

$$\mathbb{E}[\exp(\lambda X)] \leq 1 + \sum_{k=2}^{\infty} \frac{\lambda^k \mathbb{E}[X^k]}{k!} \leq 1 + \sum_{k=2}^{\infty} \frac{(k|\lambda|\sigma)^k}{k!} \leq 1 + \sum_{k=2}^{\infty} (e|\lambda|\sigma)^k,$$

where we have used that $k! \geq (k/e)^k$. When $|\lambda| < \frac{1}{e\sigma}$, evaluating the geometric series yields

$$\mathbb{E}[\exp(\lambda X)] \leq 1 + \frac{(e\lambda\sigma)^2}{1 - e|\lambda|\sigma}.$$

For $|\lambda| \leq \frac{1}{2e\sigma}$, we obtain $\mathbb{E}[e^{\lambda X}] \leq 1 + 2e^2\sigma^2\lambda^2$, and as $1 + x \leq e^x$ this implies (4).

For the opposite direction, assume (4) holds with $K_4 = K'_4 = 1$. Then $\mathbb{E}[\exp(\lambda X/\sigma)] \leq \exp(1)$ for $\lambda \in [-1, 1]$, and (3) holds. The preceding parts imply the remainder of the equivalence.

4.4.3 Proof of Theorem 5.1.6

JCD Comment: I would like to write this. For now, check out Ledoux and Talagrand [135, Theorem 4.12] or Koltchinskii [127, Theorem 2.2].

4.4.4 Proof of Proposition 4.3.2

The key insight is to rewrite the matrix exponential e^{A+B} as a limit of sums of matrices, then work more directly with traces of powers. To that end, we shall use the Lie product formula

$$\lim_{n \rightarrow \infty} (\exp(A/n) \exp(B/n))^n = \exp(A + B). \quad (4.4.1)$$

We leave the proof of the equality (4.4.1) as Exercise 4.15. Using it, however, it is evidently sufficient to prove that there exists some sequence of integers $n \rightarrow \infty$ where along this sequence,

$$\text{tr} \left(\left(e^{A/n} e^{B/n} \right)^n \right) \leq \text{tr}(e^A e^B). \quad (4.4.2)$$

Now recall that the Schatten p -norm of a matrix A is $\|A\|_p := \text{tr}((AA^*)^{p/2})^{1/p} = \|\gamma(A)\|_p$, the ℓ_p -norm of its singular values, where $p = 2$ gives the Euclidean or Frobenius norm $\|A\|_2 = (\sum_{i,j} |A_{ij}|^2)^{1/2}$. This norm gives a generalized Hölder-type inequality for powers of 2, that is, $n \in \{2^k\}_{k \in \mathbb{N}}$, which we can in turn use to prove the Golden-Thompson inequality. In particular, we demonstrate that for n a power of 2,

$$|\text{tr}(A_1 \cdots A_n)| \leq \|A_1\|_n \cdots \|A_n\|_n. \quad (4.4.3)$$

To see this inequality, we proceed inductively. Because the trace defines the inner product $\langle A, B \rangle = \text{tr}(A^*B)$, for $n = 2$, the Cauchy-Schwarz inequality implies

$$|\text{tr}(A_1 A_2)| = |\langle A_1^*, A_2 \rangle| \leq \|A_1\|_2 \|A_2\|_2.$$

We now perform an induction, where we have demonstrated the base case $n = 2$. Then for $n \geq 4$ a power of 2, we have by the inductive hypothesis that inequality (4.4.3) holds for $n/2$ that

$$|\text{tr}(A_1 \cdots A_n)| \leq \|A_1 A_2\|_{n/2} \cdots \|A_{n-1} A_n\|_{n/2}.$$

Now consider an arbitrary pair of matrices A, B . We will demonstrate that $\|AB\|_{n/2} \leq \|A\|_n \|B\|_n$, which will then evidently imply inequality (4.4.3). For these, we have

$$\|AB\|_{n/2}^{n/2} = \text{tr}(\underbrace{ABB^*A^* \cdots ABB^*A^*}_{n/4 \text{ times}}) = \text{tr}\left((A^*ABB^*)^{n/4}\right)$$

by the cyclic property of the trace. Using the inductive hypothesis again with $n/4$ copies of each of the matrices $A^T A$ and BB^T , we thus have

$$\begin{aligned} \|AB\|_{n/2}^{n/2} &\leq \text{tr}\left((A^*ABB^*)^{n/4}\right) \\ &\leq \|A^*A\|_{n/2}^{n/4} \|BB^*\|_{n/2}^{n/4} = \text{tr}\left((A^*A)^{n/2}\right)^{1/2} \text{tr}\left((BB^*)^{n/2}\right)^{1/2} = \|A\|_n^{n/2} \|B\|_n^{n/2}. \end{aligned}$$

That is, we have $\|AB\|_{n/2} \leq \|A\|_n \|B\|_n$ for any A, B as desired, giving inequality (4.4.3).

We apply inequality (4.4.3) to powers of products of Hermitian matrices A, B . We have

$$\text{tr}((AB)^n) \leq \|AB\|_n^n = \text{tr}\left((ABB^*A^*)^{n/2}\right) = \text{tr}\left((A^*ABB^*)^{n/2}\right) = \text{tr}\left((A^2B^2)^{n/2}\right)$$

because $A = A^*$ and $B = B^*$. Recognizing that A^2 and B^2 are Hermitian, we repeat this argument to obtain

$$\text{tr}\left((A^2B^2)^{n/2}\right) \leq \text{tr}\left((A^4B^4)^{n/4}\right) \leq \cdots \leq \text{tr}(A^n B^n)$$

for any $n \in \{2^k\}_{k \in \mathbb{N}}$. Replacing A and B by e^A and e^B , which are both symmetric, we obtain

$$\text{tr}((e^A e^B)^n) \leq \text{tr}(e^{nA} e^{nB}) \quad \text{for } n \in \{2^k\}_{k \in \mathbb{N}}.$$

This is inequality (4.4.2) once we replace A and B by A/n and B/n .

4.5 Bibliography

A few references on concentration, random matrices, and entropies include Vershynin's extraordinarily readable lecture notes [187], upon which our proof of Theorem 4.1.11 is based, the comprehensive book of Boucheron, Lugosi, and Massart [37], and the more advanced material in Buldygin and Kozachenko [43]. Many of our arguments are based off of those of Vershynin and Boucheron et al. Kolmogorov and Tikhomirov [126] introduced metric entropy.

We give weaker versions of the matrix-Hoeffding and matrix-Bernstein inequalities. It is possible to do much better.

Ahlsvede-Winter developed the matrix concentration inequalities and Petz [153].

I took the proof of Golden-Thompson from Terry Tao's blog.

Lemma 4.3.3 is [181, Lemma 7.6].

It is possible to obtain better concentration guarantees using Lieb's concavity inequality that

$$f(A) := \text{tr}(\exp(H + \log A)) \quad (4.5.1)$$

is a concave function in $A \succ 0$.

4.6 Exercises

Exercise 4.1 (Concentration of bounded random variables): Let X be a random variable taking values in $[a, b]$, where $-\infty < a \leq b < \infty$. In this question, we show *Hoeffding's Lemma*, that is, that X is sub-Gaussian: for all $\lambda \in \mathbb{R}$, we have

$$\mathbb{E}[\exp(\lambda(X - \mathbb{E}[X]))] \leq \exp\left(\frac{\lambda^2(b-a)^2}{8}\right).$$

(a) Show that $\text{Var}(X) \leq (\frac{b-a}{2})^2 = \frac{(b-a)^2}{4}$ for any random variable X taking values in $[a, b]$.

(b) Let

$$\phi(\lambda) = \log \mathbb{E}[\exp(\lambda(X - \mathbb{E}[X]))].$$

Assuming that $\mathbb{E}[X] = 0$ (convince yourself that this is no loss of generality) show that

$$\phi(0) = 0, \quad \phi'(0) = 0, \quad \phi''(t) = \frac{\mathbb{E}[X^2 e^{tX}]}{\mathbb{E}[e^{tX}]} - \frac{\mathbb{E}[X e^{tX}]^2}{\mathbb{E}[e^{tX}]^2}.$$

(You may assume that derivatives and expectations commute, which they do in this case.)

(c) Construct a random variable Y_t , defined for $t \in \mathbb{R}$, such that $Y_t \in [a, b]$ and

$$\text{Var}(Y_t) = \phi''(t).$$

(You may assume X has a density for simplicity.)

(d) Using the result of part (c), show that $\phi(\lambda) \leq \frac{\lambda^2(b-a)^2}{8}$ for all $\lambda \in \mathbb{R}$.

Exercise 4.2 (Variance lower bounds on sub-Gaussian parameters):

(a) Let X be σ^2 -sub-Gaussian, that is, $\mathbb{E}[\exp(\lambda X)] \leq \exp(\frac{\lambda^2 \sigma^2}{2})$ for all $\lambda \in \mathbb{R}$. Show that $\mathbb{E}[X] = 0$ and $\mathbb{E}[X^2] \leq \sigma^2$.

(b) Let X be a random variable with $\mathbb{E}[X] = 0$ and $\text{Var}(X) = \sigma^2 > 0$. Show that

$$\liminf_{|\lambda| \rightarrow \infty} \frac{1}{|\lambda|} \log \mathbb{E}[\exp(\lambda X)] > 0.$$

Exercise 4.3 (Mills ratio): Let $\phi(t) = \frac{1}{\sqrt{2\pi}} e^{-t^2/2}$ be the density of a standard Gaussian, $Z \sim \mathcal{N}(0, 1)$, and $\Phi(t) = \int_{-\infty}^t \phi(u) du$ its cumulative distribution function.

(a) Show that $\mathbb{P}(Z \geq t) \leq \frac{1}{t} \phi(t)$ for all $t > 0$.

(b) Define

$$g(t) := 1 - \Phi(t) - \frac{t}{t^2 + 1} \phi(t).$$

Show that $g(0) = 0$, $g'(t) < 0$ for all $t \geq 0$, and that $\lim_{t \rightarrow \infty} g(t) = 0$.

(c) Conclude that for all $t \geq 0$,

$$\frac{t}{t^2 + 1} \phi(t) \leq \mathbb{P}(Z \geq t) \leq \frac{1}{t} \phi(t).$$

Exercise 4.4 (Likelihood ratio bounds and concentration): Consider a data release problem, where given a sample x , we release a sequence of data Z_1, Z_2, \dots, Z_n belonging to a discrete set \mathcal{Z} , where Z_i may depend on Z_1^{i-1} and x . We assume that the data has limited information about x in the sense that for any two samples x, x' , we have the likelihood ratio bound

$$\frac{p(z_i \mid x, z_1^{i-1})}{p(z_i \mid x', z_1^{i-1})} \leq e^\varepsilon.$$

Let us control the amount of “information” (in the form of an updated log-likelihood ratio) released by this sequential mechanism. Fix x, x' , and define

$$L(z_1, \dots, z_n) := \log \frac{p(z_1, \dots, z_n \mid x)}{p(z_1, \dots, z_n \mid x')}.$$

(a) Show that, assuming the data Z_i are drawn conditional on x ,

$$\mathbb{P}(L(Z_1, \dots, Z_n) \geq n\varepsilon(e^\varepsilon - 1) + t) \leq \exp\left(-\frac{t^2}{2n\varepsilon^2}\right).$$

Equivalently, show that

$$\mathbb{P}\left(L(Z_1, \dots, Z_n) \geq n\varepsilon(e^\varepsilon - 1) + \varepsilon\sqrt{2n \log(1/\delta)}\right) \leq \delta.$$

(b) Let $\gamma \in (0, 1)$. Give the largest value of ε you can that is sufficient to guarantee that for any test $\Psi : \mathcal{Z}^n \rightarrow \{x, x'\}$, we have

$$P_x(\Psi(Z_1^n) \neq x) + P_{x'}(\Psi(Z_1^n) \neq x') \geq 1 - \gamma,$$

where P_x and $P_{x'}$ denote the sampling distribution of Z_1^n under x and x' , respectively?

Exercise 4.5 (Marcinkiewicz-Zygmund inequality): Let X_i be independent random variables with $\mathbb{E}[X_i] = 0$ and $\mathbb{E}[|X_i|^p] < \infty$, where $1 \leq p < \infty$. Prove that

$$\mathbb{E} \left[\left| \sum_{i=1}^n X_i \right|^p \right] \leq C_p \mathbb{E} \left[\left(\sum_{i=1}^n |X_i|^2 \right)^{p/2} \right]$$

where C_p is a constant depending only on p . As a corollary, derive that if $\mathbb{E}[|X_i|^p] \leq \sigma^p$ and $p \geq 2$, then

$$\mathbb{E} \left[\left| \frac{1}{n} \sum_{i=1}^n X_i \right|^p \right] \leq C_p \frac{\sigma^p}{n^{p/2}}.$$

That is, sample means converge quickly to zero in higher moments. *Hint:* For any fixed $x \in \mathbb{R}^n$, if ε_i are i.i.d. uniform signs $\varepsilon_i \in \{\pm 1\}$, then $\varepsilon^T x$ is sub-Gaussian.

Exercise 4.6 (A vector Marcinkiewicz-Zygmund inequality): Let $X_i \in \mathbb{R}^d$ be independent vectors with $\mathbb{E}[X_i] = 0$ and $\mathbb{E}[\|X_i\|_2^p] < \infty$, where $1 \leq p < \infty$. Prove that

$$\mathbb{E} \left[\left\| \sum_{i=1}^n X_i \right\|_2^p \right] \leq C_p \mathbb{E} \left[\left(\sum_{i=1}^n \|X_i\|_2^2 \right)^{p/2} \right]$$

where C_p is a constant depending only on p .

Exercise 4.7 (Small balls and anti-concentration): Let X be a nonnegative random variable satisfying $\mathbb{P}(X \leq \epsilon) \leq c\epsilon$ for some $c < \infty$ and all $\epsilon > 0$. Argue that if X_i are i.i.d. copies of X , then

$$\mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n X_i \geq t \right) \geq 1 - \exp(-2n [1/2 - 2ct]_+^2)$$

for all t .

Exercise 4.8 (Lipschitz functions remain sub-Gaussian): Let X be σ^2 -sub-Gaussian and $f : \mathbb{R} \rightarrow \mathbb{R}$ be L -Lipschitz, meaning that $|f(x) - f(y)| \leq L|x - y|$ for all x, y . Prove that there exists a numerical constant $C < \infty$ such that $f(X)$ is $CL^2\sigma^2$ -sub-Gaussian.

Exercise 4.9 (Sub-gaussian maxima): Let X_1, \dots, X_n be σ^2 -sub-gaussian (not necessarily independent) random variables. Show that

(a) $\mathbb{E}[\max_i X_i] \leq \sqrt{2\sigma^2 \log n}$.

(b) There exists a numerical constant $C < \infty$ such that $\mathbb{E}[\max_i |X_i|^p] \leq (Cp\sigma^2 \log k)^{p/2}$.

Exercise 4.10: Let $Z \sim \mathcal{N}(0, 1)$.

(a) Use the Cauchy-Schwarz inequality to show that

$$\mathbb{E} [\exp(\lambda((\mu + Z)^2 - \mu^2))] \leq \exp \left(4\lambda^2\mu^2 + \lambda + \frac{2\lambda^2}{[1 - 4|\lambda|]_+} \right).$$

(b) Use a direct integration argument, as in Examples 4.1.12 and 4.1.13, to show that

$$\mathbb{E} [\exp(\lambda(\mu + Z)^2)] \leq \exp \left(\lambda\mu^2 + \frac{2\lambda^2}{1 - 2\lambda} - \frac{1}{2} \log(1 - 2\lambda) \right)$$

for $\lambda < \frac{1}{2}$. Use this to prove the first part of Corollary 4.1.19.

Hint. It may be useful to use that $-\log(1-x) \leq -x + \frac{x^2}{2|1-x|_+}$ for all $x \in \mathbb{R}$.

Exercise 4.11: Let $Z \sim \mathcal{N}(0, I_d)$, and let $A \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$ be otherwise arbitrary. Using the first part of Corollary 4.1.19, show the second part of Corollary 4.1.19, that is, that $\|AZ - b\|_2^2$ is $(4(\|A\|_{\text{Fr}}^2 + 2\|b\|_2^2), 4\|A\|_{\text{op}}^2)$ -sub-exponential. *Hint.* Use the singular value decomposition $A = UTV^\top$ of A , and note that $V^\top Z \sim \mathcal{N}(0, I_d)$. Then

Exercise 4.12 (Sub-Gaussian constants of Bernoulli random variables): In this exercise, we will derive sharp sub-Gaussian constants for Bernoulli random variables (cf. [112, Thm. 1] or [125, 24]), showing

$$\log \mathbb{E}[e^{t(X-p)}] \leq \frac{1-2p}{4 \log \frac{1-p}{p}} t^2 \quad \text{for all } t \geq 0. \quad (4.6.1)$$

(a) Define $\varphi(t) = \log(\mathbb{E}[e^{t(X-p)}]) = \log((1-p)e^{-tp} + pe^{t(1-p)})$. Show that

$$\varphi'(t) = \mathbb{E}[Y_t] \quad \text{and} \quad \varphi''(t) = \text{Var}(Y_t)$$

where $Y_t = (1-p)$ with probability $q(t) := \frac{pe^{t(1-p)}}{pe^{t(1-p)} + (1-p)e^{-tp}}$ and $Y_t = -p$ otherwise.

(b) Show that $\varphi'(0) = 0$ and that if $p > \frac{1}{2}$, then $\text{Var}(Y_t) \leq \text{Var}(Y_0) = p(1-p)$. Conclude that $\varphi(t) \leq \frac{p(1-p)}{2} t^2$ for all $t \geq 0$.

(c) Argue that $p(1-p) \leq \frac{1-2p}{2 \log \frac{1-p}{p}}$ for $p \in [0, 1]$. *Hint:* Let $p = \frac{1+\delta}{2}$ for $\delta \in [0, 1]$, so that the inequality is equivalent to $\log \frac{1+\delta}{1-\delta} \leq \frac{2\delta}{1-\delta^2}$. Then use that $\log(1+\delta) = \int_0^\delta \frac{1}{1+u} du$.

(d) Let $C = 2 \log \frac{1-p}{p}$ and define $s = Ct = 2 \log \frac{1-p}{p} s$, and let

$$f(s) = \frac{1-2p}{2} Cs^2 + Cps - \log(1-p + pe^{Cs}),$$

so that inequality (4.6.1) holds if and only if $f(s) \geq 0$ for all $s \geq 0$. Give $f'(s)$ and $f''(s)$.

(e) Show that $f(0) = f(1) = f'(0) = f'(1) = 0$, and argue that $f''(s)$ changes signs at most twice and that $f''(0) = f''(1) > 0$. Use this to show that $f(s) \geq 0$ for all $s \geq 0$.

JCD Comment: Perhaps use transportation inequalities to prove this bound, and also maybe give Ordentlich and Weinberger's "A Distribution Dependent Refinement of Pinsker's Inequality" as an exercise.

Exercise 4.13: Let $s(p) = \frac{1-2p}{\log \frac{1-p}{p}}$. Show that s is concave on $[0, 1]$.

Exercise 4.14 (Inner products on complex matrices): Recall that $\langle \cdot, \cdot \rangle$ is a *complex* inner product on a vector space V if it satisfies the following for all $x, y, z \in V$:

(i) $\langle x, x \rangle \geq 0$, with $\langle x, x \rangle = 0$ if and only if $x = 0$.

(ii) It is conjugate symmetric, so that $\langle x, y \rangle = \overline{\langle y, x \rangle}$.

(iii) It is conjugate linear in its first argument, so that $\langle \alpha x + y, z \rangle = \overline{\alpha} \langle x, z \rangle + \langle y, z \rangle$ for all $\alpha \in \mathbb{C}$.

The vector space V is real with real inner product if property (ii) is replaced with the symmetry $\langle x, y \rangle = \langle y, x \rangle$ and linearity (iii) holds for $\alpha \in \mathbb{R}$.

- (a) Show that the space of complex $m \times n$ matrices $\mathbb{C}^{m \times n}$ has complex inner product $\langle A, B \rangle := \text{tr}(A^* B)$.
- (b) Show that the space \mathcal{H}_n of $n \times n$ Hermitian matrices is a real vector space with inner product $\langle A, B \rangle := \text{tr}(A^* B)$, and that consequently $\langle A, B \rangle \in \mathbb{R}$.

Exercise 4.15 (The Lie product formula): Let A and B be symmetric (or Hermitian) matrices.

- (a) Prove the Lie product formula (4.4.1), that is,

$$\lim_{n \rightarrow \infty} (\exp(A/n) \exp(B/n))^n = \exp(A) \exp(B).$$

Hint. One argument proceeds as follows. Let $O(\epsilon)$ denote a matrix E such that $\|E\|_{\text{op}} \lesssim \epsilon$. First, demonstrate that

$$e^{A/n} = I + \frac{1}{n}A + O(n^{-2}).$$

Then show that for any matrix A , we have $(I + n^{-1}A + o(n^{-1}))^n \rightarrow \exp(A)$. Combine these.

- (b) Give an example of matrices A and B that do not commute and for which $\exp(A + B) \neq \exp(A) \exp(B)$.

Exercise 4.16: Define the trace exponential function $f(X) := \text{tr}(e^X)$ on the Hermitian matrices.

- (a) Prove that f is monotone for the semidefinite order, that is, if $X \preceq Y$, then $f(X) \leq f(Y)$.
Hint. It is enough to show that for any $A \succeq 0$, the one-dimensional function $h(t) := f(X + tA)$ is monotone in t , or even that $h'(0) \geq 0$.
- (b) Prove that the trace exponential is convex on the Hermitian matrices, that is, $f(X) := \text{tr}(e^X)$ is convex. *Hint.* It is enough to show that for any X, V Hermitian that $h(t) := f(X + tV)$ is convex in t , for which it in turn suffices to show that $h''(0) \geq 0$.

Exercise 4.17 (The matrix-Bernstein inequality): In this question, we prove Theorem 4.3.6.

- (a) Let $X_i \in \mathcal{H}_d$ be independent Hermitian matrices and $S_n = \sum_{i=1}^n X_i$. Use the Golden-Thompson inequality (Proposition 4.3.2) to show that for all $\lambda \in \mathbb{R}$,

$$\text{tr}(\mathbb{E}[e^{\lambda S_n}]) \leq d \prod_{i=1}^n \left\| \mathbb{E}[e^{\lambda X_i}] \right\|_{\text{op}}.$$

- (b) Extend Example 4.1.14 to the matrix-valued case. Demonstrate that if X is a mean-zero Hermitian random matrix with $\|X\|_{\text{op}} \leq b$, then for all $|\lambda| < \frac{3}{b}$,

$$\mathbb{E}[\exp(\lambda X)] \preceq I + \frac{1}{1 - b|\lambda|/3} \frac{\lambda^2 \mathbb{E}[X^2]}{2}.$$

(c) Use parts (a) and (b) to show that

$$\mathrm{tr}(\mathbb{E}[e^{\lambda S_n}]) \leq d \exp \left(\lambda^2 \sum_{i=1}^n \sigma_i^2 \right)$$

for $|\lambda| \leq \frac{3}{2b}$.

(d) Prove Theorem 4.3.6.

Exercise 4.18: In this question, we use Lieb's concavity inequality (4.5.1) to obtain a stronger matrix Hoeffding inequality. Let $X_1, \dots, X_n \in \mathcal{H}_d$ be an independent sequence of $d \times d$ mean-zero Hermitian matrices. Let $S_n = \sum_{i=1}^n X_i$ be their sum.

(a) Let X be a random Hermitian matrix and $X^2 \preceq A^2$. Show that

$$\log \mathbb{E}[e^{\lambda X} \mid X] \preceq \frac{\lambda^2}{2} A^2.$$

Hint. Use that the matrix logarithm is operator monotone, that is, if $A \preceq B$, then $\log A \preceq \log B$.

(b) Assume that $X_n^2 \preceq A_n^2$. Show that

$$\mathbb{E}[\mathrm{tr}(\exp(\lambda S_n))] \leq \mathbb{E}[\mathrm{tr} \exp(\lambda S_{n-1} + 2\lambda^2 A_n^2)].$$

(c) Show that if $X_i^2 \preceq A_i^2$ for each i , then

$$\mathbb{E}[\mathrm{tr}(\exp(\lambda S_n))] \leq \exp \left(2\lambda^2 \sum_{i=1}^n A_i^2 \right).$$

(d) Show that if $\sigma^2 \geq \|\sum_{i=1}^n A_i^2\|_{\mathrm{op}}$, then for $t \geq 0$,

$$\mathbb{P}(\lambda_{\max}(S_n) \geq t) \leq d \exp \left(-\frac{t^2}{8\sigma^2} \right).$$

(e) Give an example of random Hermitian matrices where the preceding bound is much sharper than Corollary 4.3.5.

Exercise 4.19: In this question, we use Lieb's concavity inequality (4.5.1) to demonstrate a sharper matrix Bernstein inequality than Theorem 4.3.6.

(a) Define the matrix cumulant generating function $\phi_{X_i}(\lambda) := \log \mathbb{E}[\exp(\lambda X_i)]$. Show that

$$\mathbb{E}[\mathrm{tr}(\exp(\lambda S_n))] \leq \mathrm{tr} \exp \left(\sum_{i=1}^n \phi_{X_i}(\lambda) \right).$$

(b) Using Exercise 4.17 part (b), show that if $\mathbb{E}[X_i] = 0$, $\mathbb{E}[X_i^2] \preceq \Sigma_i$, and $\|X_i\|_{\mathrm{op}} \leq b$ for each i , then for $|\lambda| < \frac{3}{b}$, we have

$$\mathbb{E}[\mathrm{tr}(e^{\lambda S_n})] \leq d \exp \left(\frac{1}{1 - b|\lambda|/3} \frac{\lambda^2}{2} \sigma^2 \right) \quad \text{where} \quad \sigma^2 := \left\| \sum_{i=1}^n \Sigma_i \right\|_{\mathrm{op}}.$$

(c) Show that there exists a numerical constant $c > 0$ such that for all $t \geq 0$,

$$\mathbb{P}\left(\|S_n\|_{\text{op}} \geq t\right) \leq 2d \exp\left(-c \min\left\{\frac{t^2}{\sigma^2}, \frac{t}{b}\right\}\right).$$

Why is this sharper than Theorem [4.3.6](#)?

Chapter 5

Estimation and generalization

5.1 Uniformity and metric entropy

Now that we have explored a variety of concentration inequalities, we show how to put them to use in demonstrating that a variety of estimation, learning, and other types of procedures have nice convergence properties. We first give a somewhat general collection of results, then delve deeper by focusing on some standard tasks from machine learning.

5.1.1 Symmetrization and uniform laws

The first set of results we consider are *uniform laws of large numbers*, where the goal is to bound means uniformly over different classes of functions. Frequently, such results are called *Glivenko-Cantelli* laws, after the original Glivenko-Cantelli theorem, which shows that empirical distributions uniformly converge. We revisit these ideas in the next chapter, where we present a number of more advanced techniques based on ideas of metric entropy (or volume-like considerations); here we present the basic ideas using our stability and bounded differencing tools.

The starting point is to define what we mean by a uniform law of large numbers. To do so, we adopt notation (as in Example 4.2.8) we will use throughout the remainder of the book, reminding readers as we go. For a sample X_1, \dots, X_n on a space \mathcal{X} , we let

$$P_n := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i}$$

denote the empirical distribution on $\{X_i\}_{i=1}^n$, where $\mathbf{1}_{X_i}$ denotes the point mass at X_i . Then for functions $f : \mathcal{X} \rightarrow \mathbb{R}$ (or more generally, any function f defined on \mathcal{X}), we let

$$P_n f := \mathbb{E}_{P_n}[f(X)] = \frac{1}{n} \sum_{i=1}^n f(X_i)$$

denote the empirical expectation of f evaluated on the sample, and we also let

$$P f := \mathbb{E}_P[f(X)] = \int f(x) dP(x)$$

denote general expectations under a measure P . With this notation, we study *uniform laws of large numbers*, which consist of proving results of the form

$$\sup_{f \in \mathcal{F}} |P_n f - P f| \rightarrow 0, \tag{5.1.1}$$

where convergence is in probability, expectation, almost surely, or with rates of convergence. When we view P_n and P as (infinite-dimensional) vectors on the space of maps from $\mathcal{F} \rightarrow \mathbb{R}$, then we may define the (semi)norm $\|\cdot\|_{\mathcal{F}}$ for any $L : \mathcal{F} \rightarrow \mathbb{R}$ by

$$\|L\|_{\mathcal{F}} := \sup_{f \in \mathcal{F}} |L(f)|,$$

in which case Eq. (5.1.1) is equivalent to proving

$$\|P_n - P\|_{\mathcal{F}} \rightarrow 0.$$

Thus, roughly, we are simply asking questions about when random vectors converge to their expectations.¹

The starting point of this investigation considers bounded random functions, that is, \mathcal{F} consists of functions $f : \mathcal{X} \rightarrow [a, b]$ for some $-\infty < a \leq b < \infty$. In this case, the bounded differences inequality (Proposition 4.2.5) immediately implies that expectations of $\|P_n - P\|_{\mathcal{F}}$ provide strong guarantees on concentration of $\|P_n - P\|_{\mathcal{F}}$.

Proposition 5.1.1. *Let \mathcal{F} be as above. Then*

$$\mathbb{P}(\|P_n - P\|_{\mathcal{F}} \geq \mathbb{E}[\|P_n - P\|_{\mathcal{F}}] + t) \leq \exp\left(-\frac{2nt^2}{(b-a)^2}\right) \quad \text{for } t \geq 0.$$

Proof Let P_n and P'_n be two empirical distributions, differing only in observation i (with X_i and X'_i). We observe that

$$\begin{aligned} \sup_{f \in \mathcal{F}} |P_n f - P f| - \sup_{f \in \mathcal{F}} |P'_n f - P f| &\leq \sup_{f \in \mathcal{F}} \{|P_n f - P f| - |P'_n f - P f|\} \\ &\leq \frac{1}{n} \sup_{f \in \mathcal{F}} |f(X_i) - f(X'_i)| \leq \frac{b-a}{n} \end{aligned}$$

by the triangle inequality. An entirely parallel argument gives the converse lower bound of $-\frac{b-a}{n}$, and thus Proposition 4.2.5 gives the result. \square

Proposition 5.1.1 shows that, to provide control over high-probability concentration of $\|P_n - P\|_{\mathcal{F}}$, it is (at least in cases where \mathcal{F} is bounded) sufficient to control the expectation $\mathbb{E}[\|P_n - P\|_{\mathcal{F}}]$. We take this approach through the remainder of this section, developing tools to simplify bounding this quantity.

Our starting points consist of a few inequalities relating expectations to *symmetrized* quantities, which are frequently easier to control than their non-symmetrized parts. This symmetrization technique is widely used in probability theory, theoretical statistics, and machine learning. The key is that for centered random variables, symmetrized quantities have, to within numerical constants, similar expectations to their non-symmetrized counterparts. Thus, in many cases, it is equivalent to analyze the symmetrized quantity and the initial quantity.

Proposition 5.1.2. *Let X_i be independent random vectors on a (Banach) space with norm $\|\cdot\|$ and let $\varepsilon_i \{-1, 1\}$ be independent random signs. Then for any $p \geq 1$,*

$$2^{-p} \mathbb{E} \left[\left\| \sum_{i=1}^n \varepsilon_i (X_i - \mathbb{E}[X_i]) \right\|^p \right] \leq \mathbb{E} \left[\left\| \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right\|^p \right] \leq 2^p \mathbb{E} \left[\left\| \sum_{i=1}^n \varepsilon_i X_i \right\|^p \right]$$

¹Some readers may worry about measurability issues here. All of our applications will be in separable spaces, so that we may take suprema with abandon without worrying about measurability, and consequently we ignore this from now on.

In the proof of the upper bound, we could also show the bound

$$\mathbb{E} \left[\left\| \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right\|^p \right] \leq 2^p \mathbb{E} \left[\left\| \sum_{i=1}^n \varepsilon_i (X_i - \mathbb{E}[X_i]) \right\|^p \right],$$

so we may analyze whichever is more convenient.

Proof We prove the right bound first. We introduce independent copies of the X_i and use these to symmetrize the quantity. Indeed, let X'_i be an independent copy of X_i , and use Jensen's inequality and the convexity of $\|\cdot\|^p$ to observe that

$$\mathbb{E} \left[\left\| \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right\|^p \right] = \mathbb{E} \left[\left\| \sum_{i=1}^n (X_i - \mathbb{E}[X'_i]) \right\|^p \right] \leq \mathbb{E} \left[\left\| \sum_{i=1}^n (X_i - X'_i) \right\|^p \right].$$

Now, note that the distribution of $X_i - X'_i$ is symmetric, so that $X_i - X'_i \stackrel{\text{dist}}{=} \varepsilon_i (X_i - X'_i)$, and thus

$$\mathbb{E} \left[\left\| \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right\|^p \right] \leq \mathbb{E} \left[\left\| \sum_{i=1}^n \varepsilon_i (X_i - X'_i) \right\|^p \right].$$

Multiplying and dividing by 2^p , Jensen's inequality then gives

$$\begin{aligned} \mathbb{E} \left[\left\| \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right\|^p \right] &\leq 2^p \mathbb{E} \left[\left\| \frac{1}{2} \sum_{i=1}^n \varepsilon_i (X_i - X'_i) \right\|^p \right] \\ &\leq 2^{p-1} \left[\mathbb{E} \left[\left\| \sum_{i=1}^n \varepsilon_i X_i \right\|^p \right] + \mathbb{E} \left[\left\| \sum_{i=1}^n \varepsilon_i X'_i \right\|^p \right] \right] \end{aligned}$$

as desired.

For the left bound in the proposition, let $Y_i = X_i - \mathbb{E}[X_i]$ be the centered version of the random variables. We break the sum over random variables into two parts, conditional on whether $\varepsilon_i = \pm 1$, using repeated conditioning. We have

$$\begin{aligned} \mathbb{E} \left[\left\| \sum_{i=1}^n \varepsilon_i Y_i \right\|^p \right] &= \mathbb{E} \left[\left\| \sum_{i:\varepsilon_i=1} Y_i - \sum_{i:\varepsilon_i=-1} Y_i \right\|^p \right] \\ &\leq \mathbb{E} \left[2^{p-1} \mathbb{E} \left[\left\| \sum_{i:\varepsilon_i=1} Y_i \right\|^p \mid \varepsilon \right] + 2^{p-1} \mathbb{E} \left[\left\| \sum_{i:\varepsilon_i=-1} Y_i \right\|^p \mid \varepsilon \right] \right] \\ &= 2^{p-1} \mathbb{E} \left[\mathbb{E} \left[\left\| \sum_{i:\varepsilon_i=1} Y_i + \sum_{i:\varepsilon_i=-1} \mathbb{E}[Y_i] \right\|^p \mid \varepsilon \right] + \mathbb{E} \left[\left\| \sum_{i:\varepsilon_i=-1} Y_i + \sum_{i:\varepsilon_i=1} \mathbb{E}[Y_i] \right\|^p \mid \varepsilon \right] \right] \\ &\leq 2^{p-1} \mathbb{E} \left[\mathbb{E} \left[\left\| \sum_{i:\varepsilon_i=1} Y_i + \sum_{i:\varepsilon_i=-1} Y_i \right\|^p \mid \varepsilon \right] + \mathbb{E} \left[\left\| \sum_{i:\varepsilon_i=-1} Y_i + \sum_{i:\varepsilon_i=1} Y_i \right\|^p \mid \varepsilon \right] \right] \\ &= 2^p \mathbb{E} \left[\left\| \sum_{i=1}^n Y_i \right\|^p \right]. \end{aligned}$$

□

We obtain as an immediate corollary a symmetrization bound for supremum norms on function spaces. In this corollary, we use the symmetrized empirical measure

$$P_n^0 := \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{1}_{X_i}, \quad P_n^0 f = \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i).$$

The expectation of $\|P_n^0\|_{\mathcal{F}}$ is of course the Rademacher complexity (Examples 4.2.7 and 4.2.8), and we have the following corollary.

Corollary 5.1.3. *Let \mathcal{F} be a class of functions $f : \mathcal{X} \rightarrow \mathbb{R}$ and X_i be i.i.d. Then $\mathbb{E}[\|P_n - P\|_{\mathcal{F}}] \leq 2\mathbb{E}[\|P_n^0\|_{\mathcal{F}}]$.*

From Corollary 5.1.3, it is evident that by controlling the *expectation* of the symmetrized process $\mathbb{E}[\|P_n^0\|_{\mathcal{F}}]$ we can derive concentration inequalities and uniform laws of large numbers. For example, we immediately obtain that

$$\mathbb{P}(\|P_n - P\|_{\mathcal{F}} \geq 2\mathbb{E}[\|P_n^0\|_{\mathcal{F}}] + t) \leq \exp\left(-\frac{2nt^2}{(b-a)^2}\right)$$

for all $t \geq 0$ whenever \mathcal{F} consists of functions $f : \mathcal{X} \rightarrow [a, b]$.

There are numerous examples of uniform laws of large numbers, many of which reduce to developing bounds on the expectation $\mathbb{E}[\|P_n^0\|_{\mathcal{F}}]$, which is frequently possible via more advanced techniques we develop in Chapter 7. A frequent application of these symmetrization ideas is to risk minimization problems, as we discuss in the coming section; for these, it will be useful for us to develop a few analytic and calculus tools. To better match the development of these ideas, we return to the notation of Rademacher complexities, so that $R_n(\mathcal{F}) := \mathbb{E}[\|P_n^0\|_{\mathcal{F}}]$. The first is a standard result, which we state for its historical value and the simplicity of its proof.

Proposition 5.1.4 (Massart's finite class bound). *Let \mathcal{F} be any collection of functions with $f : \mathcal{X} \rightarrow \mathbb{R}$, and assume that $\sigma_n^2 := n^{-1} \mathbb{E}[\max_{f \in \mathcal{F}} \sum_{i=1}^n f(X_i)^2] < \infty$. Then*

$$R_n(\mathcal{F}) \leq \frac{\sqrt{2\sigma_n^2 \log |\mathcal{F}|}}{\sqrt{n}}.$$

Proof For each fixed x_1^n , the random variable $\sum_{i=1}^n \varepsilon_i f(x_i)$ is $\sum_{i=1}^n f(x_i)^2$ -sub-Gaussian. Now, define $\sigma^2(x_1^n) := n^{-1} \max_{f \in \mathcal{F}} \sum_{i=1}^n f(x_i)^2$. Using the results of Exercise 4.9, that is, that $\mathbb{E}[\max_{j \leq n} Z_j] \leq \sqrt{2\sigma^2 \log n}$ if the Z_j are each σ^2 -sub-Gaussian, we see that

$$R_n(\mathcal{F} \mid x_1^n) \leq \frac{\sqrt{2\sigma^2(x_1^n) \log |\mathcal{F}|}}{\sqrt{n}}.$$

Jensen's inequality that $\mathbb{E}[\sqrt{\cdot}] \leq \sqrt{\mathbb{E}[\cdot]}$ gives the result. \square

A refinement of Massart's finite class bound applies when the classes are infinite but, on a collection X_1, \dots, X_n , the functions $f \in \mathcal{F}$ may take on only a (smaller) number of values. In this case, we define the *empirical shatter coefficient* of a collection of points x_1, \dots, x_n by $S_{\mathcal{F}}(x_1^n) := \text{card}\{(f(x_1), \dots, f(x_n)) \mid f \in \mathcal{F}\}$, the number of distinct vectors of values $(f(x_1), \dots, f(x_n))$ the functions $f \in \mathcal{F}$ may take. The *shatter coefficient* is the maximum of the empirical shatter coefficients over $x_1^n \in \mathcal{X}^n$, that is, $S_{\mathcal{F}}(n) := \sup_{x_1^n} S_{\mathcal{F}}(x_1^n)$. It is clear that $S_{\mathcal{F}}(n) \leq |\mathcal{F}|$ always, but by only counting distinct values, we have the following corollary.

Corollary 5.1.5 (A sharper variant of Massart's finite class bound). *Let \mathcal{F} be any collection of functions with $f : \mathcal{X} \rightarrow \mathbb{R}$, and assume that $\sigma_n^2 := n^{-1}\mathbb{E}[\max_{f \in \mathcal{F}} \sum_{i=1}^n f(X_i)^2] < \infty$. Then*

$$R_n(\mathcal{F}) \leq \frac{\sqrt{2\sigma_n^2 \log \mathcal{S}_{\mathcal{F}}(n)}}{\sqrt{n}}.$$

Typical classes with small shatter coefficients include Vapnik-Chervonenkis classes of functions; we do not discuss these further here, instead referring to one of the many books in machine learning and empirical process theory in statistics.

The most important of the calculus rules we use are the *comparison inequalities* for Rademacher sums, which allow us to consider compositions of function classes and maintain small complexity measurers. We state the rule here; the proof is complex, so we defer it to Section 4.4.3

Theorem 5.1.6 (Ledoux-Talagrand Contraction). *Let $T \subset \mathbb{R}^n$ be an arbitrary set and let $\phi_i : \mathbb{R} \rightarrow \mathbb{R}$ be 1-Lipschitz and satisfy $\phi_i(0) = 0$. Then for any nondecreasing convex function $\Phi : \mathbb{R} \rightarrow \mathbb{R}_+$,*

$$\mathbb{E} \left[\Phi \left(\frac{1}{2} \sup_{t \in T} \left| \sum_{i=1}^n \phi_i(t_i) \varepsilon_i \right| \right) \right] \leq \mathbb{E} \left[\Phi \left(\sup_{t \in T} \langle t, \varepsilon \rangle \right) \right].$$

A corollary to this theorem is suggestive of its power and applicability. Let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ be L -Lipschitz, and for a function class \mathcal{F} define $\phi \circ \mathcal{F} = \{\phi \circ f \mid f \in \mathcal{F}\}$. Then we have the following corollary about Rademacher complexities of contractive mappings.

Corollary 5.1.7. *Let \mathcal{F} be an arbitrary function class and ϕ be L -Lipschitz. Then*

$$R_n(\phi \circ \mathcal{F}) \leq 2LR_n(\mathcal{F}) + |\phi(0)|/\sqrt{n}.$$

Proof The result is an almost immediate consequence of Theorem 5.1.6; we simply recenter our functions. Indeed, we have

$$\begin{aligned} R_n(\phi \circ \mathcal{F} \mid x_1^n) &= \mathbb{E} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (\phi(f(x_i)) - \phi(0)) + \frac{1}{n} \sum_{i=1}^n \varepsilon_i \phi(0) \right| \right] \\ &\leq \mathbb{E} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (\phi(f(x_i)) - \phi(0)) \right| \right] + \mathbb{E} \left[\left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \phi(0) \right| \right] \\ &\leq 2LR_n(\mathcal{F}) + \frac{|\phi(0)|}{\sqrt{n}}, \end{aligned}$$

where the final inequality follows by Theorem 5.1.6 (as $g(\cdot) = \phi(\cdot) - \phi(0)$ is Lipschitz and satisfies $g(0) = 0$) and that $\mathbb{E}[|\sum_{i=1}^n \varepsilon_i|] \leq \sqrt{n}$. \square

5.1.2 Metric entropy, coverings, and packings

When the class of functions \mathcal{F} under consideration is finite, the union bound more or less provides guarantees that $P_n f$ is uniformly close to $P f$ for all $f \in \mathcal{F}$. When \mathcal{F} is infinite, however, we require a different set of tools for addressing uniform laws. In many cases, because of the application of the bounded differences inequality in Proposition 5.1.1, all we really need to do is to control

the expectation $\mathbb{E}[\|P_n^0\|_{\mathcal{F}}]$, though the techniques we develop here will have broader use and can sometimes directly guarantee concentration.

The basic object we wish to control is a measure of the size of the space on which we work. To that end, we modify notation a bit to simply consider arbitrary vectors $\theta \in \Theta$, where Θ is a non-empty set with an associated (semi)metric ρ . For many purposes in estimation (and in our optimality results in the further parts of the book), a natural way to measure the size of the set is via the number of balls of a fixed radius $\delta > 0$ required to cover it.

Definition 5.1 (Covering number). *Let Θ be a set with (semi)metric ρ . A δ -cover of the set Θ with respect to ρ is a set $\{\theta_1, \dots, \theta_N\}$ such that for any point $\theta \in \Theta$, there exists some $v \in \{1, \dots, N\}$ such that $\rho(\theta, \theta_v) \leq \delta$. The δ -covering number of Θ is*

$$N(\delta, \Theta, \rho) := \inf \{N \in \mathbb{N} : \text{there exists a } \delta\text{-cover } \theta_1, \dots, \theta_N \text{ of } \Theta\}.$$

The *metric entropy* of the set Θ is simply the logarithm of its covering number $\log N(\delta, \Theta, \rho)$. We can define a related measure—more useful for constructing our lower bounds—of size that relates to the number of disjoint balls of radius $\delta > 0$ that can be placed into the set Θ .

Definition 5.2 (Packing number). *A δ -packing of the set Θ with respect to ρ is a set $\{\theta_1, \dots, \theta_M\}$ such that for all distinct $v, v' \in \{1, \dots, M\}$, we have $\rho(\theta_v, \theta_{v'}) \geq \delta$. The δ -packing number of Θ is*

$$M(\delta, \Theta, \rho) := \sup \{M \in \mathbb{N} : \text{there exists a } \delta\text{-packing } \theta_1, \dots, \theta_M \text{ of } \Theta\}.$$

Figures 5.1 and 5.2 give examples of (respectively) a covering and a packing of the same set.

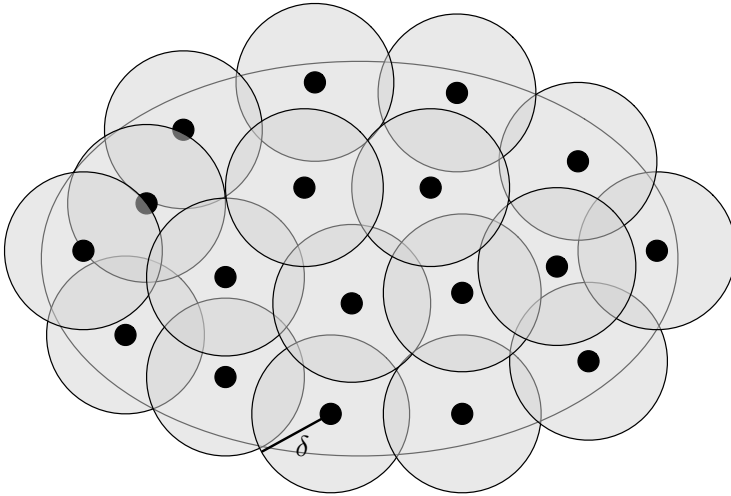


Figure 5.1. A δ -covering of the elliptical set by balls of radius δ .

An exercise in proof by contradiction shows that the packing and covering numbers of a set are in fact closely related:

Lemma 5.1.8. *The packing and covering numbers satisfy the following inequalities:*

$$M(2\delta, \Theta, \rho) \leq N(\delta, \Theta, \rho) \leq M(\delta, \Theta, \rho).$$

We leave derivation of this lemma to Exercise 5.2, noting that it shows that (up to constant factors) packing and covering numbers have the same scaling in the radius δ . As a simple example, we see

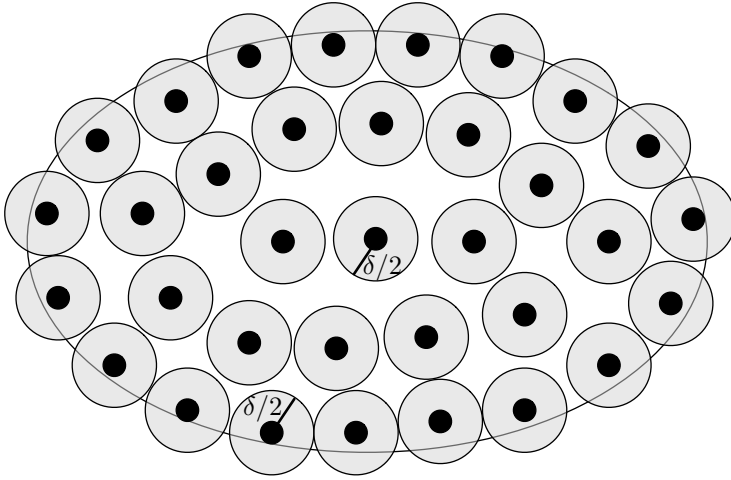


Figure 5.2. A δ -packing of the elliptical set, where balls have radius $\delta/2$. No balls overlap, and each center of the packing satisfies $\|\theta_v - \theta_{v'}\| \geq \delta$.

for any interval $[a, b]$ on the real line that in the usual absolute distance metric, $N(\delta, [a, b], |\cdot|) \asymp (b - a)/\delta$.

As one example of the metric entropy, consider a set of functions \mathcal{F} with reasonable covering numbers (metric entropy) in $\|\cdot\|_\infty$ -norm.

Example 5.1.9 (The “standard” covering number guarantee): Let \mathcal{F} consist of functions $f : \mathcal{X} \rightarrow [-b, b]$ and let the metric ρ be $\|f - g\|_\infty = \sup_{x \in \mathcal{X}} |f(x) - g(x)|$. Then

$$\mathbb{P} \left(\sup_{f \in \mathcal{F}} |P_n f - P f| \geq t \right) \leq \exp \left(-\frac{nt^2}{18b^2} + \log N(t/3, \mathcal{F}, \|\cdot\|_\infty) \right). \quad (5.1.2)$$

So as long as the covering numbers $N(t, \mathcal{F}, \|\cdot\|_\infty)$ grow sub-exponentially in t —so that $\log N(t) \ll nt^2$ —we have the (essentially) sub-Gaussian tail bound (5.1.2). Example 5.2.11 gives one typical case. Indeed, fix a minimal $t/3$ -cover of \mathcal{F} in $\|\cdot\|_\infty$ of size $N := N(t/3, \mathcal{F}, \|\cdot\|_\infty)$, calling the covering functions f_1, \dots, f_N . Then for any $f \in \mathcal{F}$ and the function f_i satisfying $\|f - f_i\|_\infty \leq t/3$, we have

$$|P_n f - P f| \leq |P_n f - P_n f_i| + |P_n f_i - P f_i| + |P f_i - P f| \leq |P_n f_i - P f_i| + \frac{2t}{3}.$$

The Azuma-Hoeffding inequality (Theorem 4.2.3) guarantees (by a union bound) that

$$\mathbb{P} \left(\max_{i \leq N} |P_n f_i - P f_i| \geq t \right) \leq \exp \left(-\frac{nt^2}{2b^2} + \log N \right).$$

Combine this bound (replacing t with $t/3$) to obtain inequality (5.1.2). \diamond

Given the relationships between packing, covering, and size of sets Θ , we would expect there to be relationships between volume, packing, and covering numbers. This is indeed the case, as we now demonstrate for arbitrary norm balls in finite dimensions.

Lemma 5.1.10. Let \mathbb{B} denote the unit $\|\cdot\|$ -ball in \mathbb{R}^d . Then

$$\left(\frac{1}{\delta} \right)^d \leq N(\delta, \mathbb{B}, \|\cdot\|) \leq \left(1 + \frac{2}{\delta} \right)^d.$$

Proof We prove the lemma via a volumetric argument. For the lower bound, note that if the points v_1, \dots, v_N are a δ -cover of \mathbb{B} , then

$$\text{Vol}(\mathbb{B}) \leq \sum_{i=1}^N \text{Vol}(\delta\mathbb{B} + v_i) = N \text{Vol}(\delta\mathbb{B}) = N \text{Vol}(\mathbb{B})\delta^d.$$

In particular, $N \geq \delta^{-d}$. For the upper bound on $N(\delta, \mathbb{B}, \|\cdot\|)$, let \mathcal{V} be a δ -packing of \mathbb{B} with maximal cardinality, so that $|\mathcal{V}| = M(\delta, \mathbb{B}, \|\cdot\|) \geq N(\delta, \mathbb{B}, \|\cdot\|)$ (recall Lemma 5.1.8). Notably, the collection of δ -balls $\{\delta\mathbb{B} + v_i\}_{i=1}^M$ cover the ball \mathbb{B} (as otherwise, we could put an additional element in the packing \mathcal{V}), and moreover, the balls $\{\frac{\delta}{2}\mathbb{B} + v_i\}$ are all disjoint by definition of a packing. Consequently, we find that

$$M \left(\frac{\delta}{2}\right)^d \text{Vol}(\mathbb{B}) = M \text{Vol}\left(\frac{\delta}{2}\mathbb{B}\right) \leq \text{Vol}\left(\mathbb{B} + \frac{\delta}{2}\mathbb{B}\right) = \left(1 + \frac{\delta}{2}\right)^d \text{Vol}(\mathbb{B}).$$

Rewriting, we obtain

$$M(\delta, \mathbb{B}, \|\cdot\|) \leq \left(\frac{2}{\delta}\right)^d \left(1 + \frac{\delta}{2}\right)^d \frac{\text{Vol}(\mathbb{B})}{\text{Vol}(\mathbb{B})} = \left(1 + \frac{2}{\delta}\right)^d,$$

completing the proof. \square

5.1.3 Application: matrix concentration

Let us give one application of Lemma 5.1.10 to concentration of random matrices; we explore more in the exercises as well. We can generalize the definition of sub-Gaussian random variables to *sub-Gaussian random vectors*, where we say that $X \in \mathbb{R}^d$ is a σ^2 -sub-Gaussian vector if

$$\mathbb{E}[\exp(\langle u, X - \mathbb{E}[X] \rangle)] \leq \exp\left(\frac{\sigma^2}{2} \|u\|_2^2\right) \quad (5.1.3)$$

for all $u \in \mathbb{R}^d$. For example, $X \sim \mathcal{N}(0, I_d)$ is immediately 1-sub-Gaussian, and $X \in [-b, b]^d$ with independent entries is b^2 -sub-Gaussian. Now, suppose that X_i are independent isotropic random vectors, meaning that $\mathbb{E}[X_i] = 0$, $\mathbb{E}[X_i X_i^\top] = I_d$, and that they are also σ^2 -sub-Gaussian. Then by an application of Lemma 5.1.10, we can give concentration guarantees for the sample covariance $\Sigma_n := \frac{1}{n} \sum_{i=1}^n X_i X_i^\top$ for the operator norm $\|A\|_{\text{op}} := \sup\{\langle u, Av \rangle \mid \|u\|_2 = \|v\|_2 = 1\}$.

Proposition 5.1.11. *Let X_i be independent isotropic and σ^2 -sub-Gaussian vectors. Then there is a numerical constant C such that the sample covariance $\Sigma_n := \frac{1}{n} \sum_{i=1}^n X_i X_i^\top$ satisfies*

$$\|\Sigma_n - I_d\|_{\text{op}} \leq C\sigma^2 \left[\frac{d + \log \frac{1}{\delta}}{n} + \sqrt{\frac{d + \log \frac{1}{\delta}}{n}} \right]$$

with probability at least $1 - \delta$.

Proof We begin with an intermediate lemma.

Lemma 5.1.12. *Let A be symmetric and $\{u_i\}_{i=1}^N$ be an ϵ -cover of the unit ℓ_2 ball \mathbb{B}_2^d . Then*

$$(1 - 2\epsilon) \|A\|_{\text{op}} \leq \max_{i \leq N} \langle u_i, Au_i \rangle \leq \|A\|_{\text{op}}.$$

Proof The second inequality is trivial. Fix any $u \in \mathbb{B}_2^d$. Then for the i such that $\|u - u_i\|_2 \leq \epsilon$, we have

$$\langle u, Au \rangle = \langle u - u_i, Au \rangle + \langle u_i, Au \rangle = 2\langle u - u_i, Au \rangle + \langle u_i, Au_i \rangle \leq 2\epsilon \|A\|_{\text{op}} + \langle u_i, Au_i \rangle$$

by definition of the operator norm. Taking a supremum over u gives the final result. \square

Let the matrix $E_i = X_i X_i^\top - I$, and define the average error $\bar{E}_n = \frac{1}{n} E_i$. Then with this lemma in hand, we see that for any ϵ -cover \mathcal{N} of the ℓ_2 -ball \mathbb{B}_2^d ,

$$(1 - 2\epsilon) \|\bar{E}_n\|_{\text{op}} \leq \max_{u \in \mathcal{N}} \langle u, \bar{E}_n u \rangle.$$

Now, note that $\langle u, E_i u \rangle = \langle u, X_i \rangle^2 - \|u\|_2^2$ is sub-exponential, as it is certainly mean 0 and, moreover, is the square of a sub-Gaussian; in particular, Theorem 4.1.15 shows that there is a numerical constant $C < \infty$ such that

$$\mathbb{E}[\exp(\lambda \langle u, E_i u \rangle)] \leq \exp(C\lambda^2 \sigma^4) \quad \text{for } |\lambda| \leq \frac{1}{C\sigma^2}.$$

Taking $\epsilon = \frac{1}{4}$ in our covering \mathcal{N} , then,

$$\mathbb{P}(\|\bar{E}_n\|_{\text{op}} \geq t) \leq \mathbb{P}\left(\max_{u \in \mathcal{N}} \langle u, \bar{E}_n u \rangle \geq t/2\right) \leq |\mathcal{N}| \cdot \max_{u \in \mathcal{N}} \mathbb{P}(\langle u, n\bar{E}_n u \rangle \geq nt/2)$$

by a union bound. As sums of sub-exponential random variable remain sub-exponential, Corollary 4.1.18 implies

$$\mathbb{P}(\|\bar{E}_n\|_{\text{op}} \geq t) \leq |\mathcal{N}| \exp\left(-c \min\left\{\frac{nt^2}{\sigma^4}, \frac{nt}{\sigma^2}\right\}\right),$$

where $c > 0$ is a numerical constant. Finally, we apply Lemma 5.1.10, which guarantees that $|\mathcal{N}| \leq 9^d$, and then take t to scale as the maximum of $\sigma^2 \frac{d + \log \frac{1}{\delta}}{n}$ and $\sigma^2 \sqrt{\frac{d + \log \frac{1}{\delta}}{n}}$. \square

5.2 Generalization bounds

We now build off of our ideas on uniform laws of large numbers and Rademacher complexities to demonstrate their applications in statistical machine learning problems, focusing on *empirical risk minimization* procedures and related problems. We consider a setting as follows: we have a sample $Z_1, \dots, Z_n \in \mathcal{Z}$ drawn i.i.d. according to some (unknown) distribution P , and we have a collection of functions \mathcal{F} from which we wish to select an f that “fits” the data well, according to some loss measure $\ell : \mathcal{F} \times \mathcal{Z} \rightarrow \mathbb{R}$. That is, we wish to find a function $f \in \mathcal{F}$ minimizing the *risk*

$$L(f) := \mathbb{E}_P[\ell(f, Z)]. \tag{5.2.1}$$

In general, however, we only have access to the risk via the empirical distribution of the Z_i , and we often choose f by minimizing the empirical risk

$$\widehat{L}_n(f) := \frac{1}{n} \sum_{i=1}^n \ell(f, Z_i). \quad (5.2.2)$$

As written, this formulation is quite abstract, so we provide a few examples to make it somewhat more concrete.

Example 5.2.1 (Binary classification problems): One standard problem—still abstract—that motivates the formulation (5.2.1) is the *binary classification problem*. Here the data Z_i come in pairs (X, Y) , where $X \in \mathcal{X}$ is some set of covariates (independent variables) and $Y \in \{-1, 1\}$ is the label of example X . The function class \mathcal{F} consists of functions $f : \mathcal{X} \rightarrow \mathbb{R}$, and the goal is to find a function f such that

$$\mathbb{P}(\text{sign}(f(X)) \neq Y)$$

is small, that is, minimizing the risk $\mathbb{E}[\ell(f, Z)]$ where the loss is the 0-1 loss, $\ell(f, (x, y)) = \mathbf{1}\{f(x)y \leq 0\}$. \diamond

Example 5.2.2 (Multiclass classification): The multiclass classification problem is identical to the binary problem, but instead of $Y \in \{-1, 1\}$ we assume that $Y \in [k] = \{1, \dots, k\}$ for some $k \geq 2$, and the function class \mathcal{F} consists of (a subset of) functions $f : \mathcal{X} \rightarrow \mathbb{R}^k$. The goal is to find a function f such that, if $Y = y$ is the correct label for a datapoint x , then $f_y(x) > f_l(x)$ for all $l \neq y$. That is, we wish to find $f \in \mathcal{F}$ minimizing

$$\mathbb{P}(\exists l \neq Y \text{ such that } f_l(X) \geq f_Y(X)).$$

In this case, the loss function is the zero-one loss $\ell(f, (x, y)) = \mathbf{1}\{\max_{l \neq y} f_l(x) \geq f_y(x)\}$. \diamond

Example 5.2.3 (Binary classification with linear functions): In the standard statistical learning setting, the data x belong to \mathbb{R}^d , and we assume that our function class \mathcal{F} is indexed by a set $\Theta \subset \mathbb{R}^d$, so that $\mathcal{F} = \{f_\theta : f_\theta(x) = \theta^\top x, \theta \in \Theta\}$. In this case, we may use the zero-one loss, the convex hinge loss, or the (convex) logistic loss, which are variously $\ell_{\text{zo}}(f_\theta, (x, y)) := \mathbf{1}\{y\theta^\top x \leq 0\}$, and the convex losses

$$\ell_{\text{hinge}}(f_\theta, (x, y)) = \left[1 - yx^\top \theta\right]_+ \quad \text{and} \quad \ell_{\text{logit}}(f_\theta, (x, y)) = \log(1 + \exp(-yx^\top \theta)).$$

The hinge and logistic losses, as they are convex, are substantially computationally easier to work with, and they are common choices in applications. \diamond

The main motivating question that we ask is the following: given a sample Z_1, \dots, Z_n , if we choose some $\widehat{f}_n \in \mathcal{F}$ based on this sample, can we guarantee that it generalizes to unseen data? In particular, can we guarantee that (with high probability) we have the empirical risk bound

$$\widehat{L}_n(\widehat{f}_n) = \frac{1}{n} \sum_{i=1}^n \ell(\widehat{f}_n, Z_i) \leq R(\widehat{f}_n) + \epsilon \quad (5.2.3)$$

for some small ϵ ? If we allow \widehat{f}_n to be arbitrary, then this becomes clearly impossible: consider the classification example 5.2.1, and set \widehat{f}_n to be the “hash” function that sets $\widehat{f}_n(x) = y$ if the pair (x, y) was in the sample, and otherwise $\widehat{f}_n(x) = -1$. Then clearly $\widehat{L}_n(\widehat{f}_n) = 0$, while there is no useful bound on $R(\widehat{f}_n)$.

5.2.1 Finite and countable classes of functions

In order to get bounds of the form (5.2.3), we require a few assumptions that are not too onerous. First, throughout this section, we will assume that for any fixed function f , the loss $\ell(f, Z)$ is σ^2 -sub-Gaussian, that is,

$$\mathbb{E}_P [\exp(\lambda(\ell(f, Z) - L(f)))] \leq \exp\left(\frac{\lambda^2 \sigma^2}{2}\right) \quad (5.2.4)$$

for all $f \in \mathcal{F}$. (Recall that the risk functional $L(f) = \mathbb{E}_P[\ell(f, Z)]$.) For example, if the loss is the zero-one loss from classification problems, inequality (5.2.4) is satisfied with $\sigma^2 = \frac{1}{4}$ by Hoeffding's lemma. In order to guarantee a bound of the form (5.2.4) for a function \hat{f} chosen dependent on the data, in this section we give uniform bounds, that is, we would like to bound

$$\mathbb{P}\left(\text{there exists } f \in \mathcal{F} \text{ s.t. } L(f) > \hat{L}_n(f) + t\right) \quad \text{or} \quad \mathbb{P}\left(\sup_{f \in \mathcal{F}} |\hat{L}_n(f) - R(f)| > t\right).$$

Such uniform bounds are certainly sufficient to guarantee that the empirical risk is a good proxy for the true risk L , even when \hat{f}_n is chosen based on the data.

Now, recalling that our set of functions or predictors \mathcal{F} is finite or countable, let us suppose that for each $f \in \mathcal{F}$, we have a complexity measure $c(f)$ —a penalty—such that

$$\sum_{f \in \mathcal{F}} e^{-c(f)} \leq 1. \quad (5.2.5)$$

This inequality should look familiar to the Kraft inequality—which we will see in the coming chapters—from coding theory. As soon as we have such a penalty function, however, we have the following result.

Theorem 5.2.4. *Let the loss ℓ , distribution P on \mathcal{Z} , and function class \mathcal{F} be such that $\ell(f, Z)$ is σ^2 -sub-Gaussian for each $f \in \mathcal{F}$, and assume that the complexity inequality (5.2.5) holds. Then with probability at least $1 - \delta$ over the sample $Z_{1:n}$,*

$$L(f) \leq \hat{L}_n(f) + \sqrt{2\sigma^2 \frac{\log \frac{1}{\delta} + c(f)}{n}} \quad \text{for all } f \in \mathcal{F}.$$

Proof First, we note that by the usual sub-Gaussian concentration inequality (Corollary 4.1.10) we have for any $t \geq 0$ and any $f \in \mathcal{F}$ that

$$\mathbb{P}\left(L(f) \geq \hat{L}_n(f) + t\right) \leq \exp\left(-\frac{nt^2}{2\sigma^2}\right).$$

Now, if we replace t by $\sqrt{t^2 + 2\sigma^2 c(f)/n}$, we obtain

$$\mathbb{P}\left(L(f) \geq \hat{L}_n(f) + \sqrt{t^2 + 2\sigma^2 c(f)/n}\right) \leq \exp\left(-\frac{nt^2}{2\sigma^2} - c(f)\right).$$

Then using a union bound, we have

$$\begin{aligned} \mathbb{P}\left(\exists f \in \mathcal{F} \text{ s.t. } L(f) \geq \hat{L}_n(f) + \sqrt{t^2 + 2\sigma^2 c(f)/n}\right) &\leq \sum_{f \in \mathcal{F}} \exp\left(-\frac{nt^2}{2\sigma^2} - c(f)\right) \\ &= \exp\left(-\frac{nt^2}{2\sigma^2}\right) \underbrace{\sum_{f \in \mathcal{F}} \exp(-c(f))}_{\leq 1}. \end{aligned}$$

Setting $t^2 = 2\sigma^2 \log \frac{1}{\delta}/n$ gives the result. \square

As one classical example of this setting, suppose that we have a finite class of functions \mathcal{F} . Then we can set $c(f) = \log |\mathcal{F}|$, in which case we clearly have the summation guarantee (5.2.5), and we obtain

$$L(f) \leq \widehat{L}_n(f) + \sqrt{2\sigma^2 \frac{\log \frac{1}{\delta} + \log |\mathcal{F}|}{n}} \quad \text{uniformly for } f \in \mathcal{F}$$

with probability at least $1 - \delta$. To make this even more concrete, consider the following example.

Example 5.2.5 (Floating point classifiers): We implement a linear binary classifier using double-precision floating point values, that is, we have $f_\theta(x) = \theta^\top x$ for all $\theta \in \mathbb{R}^d$ that may be represented using d double-precision floating point numbers. Then for each coordinate of θ , there are at most 2^{64} representable numbers; in total, we must thus have $|\mathcal{F}| \leq 2^{64d}$. Thus, for the zero-one loss $\ell_{zo}(f_\theta, (x, y)) = \mathbf{1}\{\theta^\top xy \leq 0\}$, we have

$$L(f_\theta) \leq \widehat{L}_n(f_\theta) + \sqrt{\frac{\log \frac{1}{\delta} + 45d}{2n}}$$

for all representable classifiers simultaneously, with probability at least $1 - \delta$, as the zero-one loss is $1/4$ -sub-Gaussian. (Here we have used that $64 \log 2 < 45$.) \diamond

We also note in passing that by replacing δ with $\delta/2$ in the bounds of Theorem 5.2.4, a union bound yields the following two-sided corollary.

Corollary 5.2.6. *Under the conditions of Theorem 5.2.4, we have*

$$\left| \widehat{L}_n(f) - L(f) \right| \leq \sqrt{2\sigma^2 \frac{\log \frac{2}{\delta} + c(f)}{n}} \quad \text{for all } f \in \mathcal{F}$$

with probability at least $1 - \delta$.

5.2.2 Large classes

When the collection of functions is (uncountably) infinite, it can be more challenging to obtain strong generalization bounds, though there still exist numerous tools for these ideas. The most basic, of which we will give examples, leverage covering number bounds (essentially, as in Example 5.1.9). We return in the next chapter to alternative approaches based on randomization and divergence measures, which provide guarantees with somewhat similar structure to those we present here.

Let us begin by considering a few examples, after which we provide examples showing how to derive explicit bounds using Rademacher complexities.

Example 5.2.7 (Rademacher complexity of the ℓ_2 -ball): Let $\Theta = \{\theta \in \mathbb{R}^d \mid \|\theta\|_2 \leq r\}$, and consider the class of linear functionals $\mathcal{F} := \{f_\theta(x) = \theta^\top x, \theta \in \Theta\}$. Then

$$R_n(\mathcal{F} \mid x_1^n) \leq \frac{r}{n} \sqrt{\sum_{i=1}^n \|x_i\|_2^2},$$

because we have

$$R_n(\mathcal{F} \mid x_1^n) = \frac{r}{n} \mathbb{E} \left[\left\| \sum_{i=1}^n \varepsilon_i x_i \right\|_2 \right] \leq \frac{r}{n} \sqrt{\mathbb{E} \left[\left\| \sum_{i=1}^n \varepsilon_i x_i \right\|_2^2 \right]} = \frac{r}{n} \sqrt{\sum_{i=1}^n \|x_i\|_2^2},$$

as desired. \diamond

In high-dimensional situations, it is sometimes useful to consider more restrictive function classes, for example, those indexed by vectors in an ℓ_1 -ball.

Example 5.2.8 (Rademacher complexity of the ℓ_1 -ball): In contrast to the previous example, suppose that $\Theta = \{\theta \in \mathbb{R}^d \mid \|\theta\|_1 \leq r\}$, and consider the linear class $\mathcal{F} := \{f_\theta(x) = \theta^T x, \theta \in \Theta\}$. Then

$$R_n(\mathcal{F} \mid x_1^n) = \frac{r}{n} \mathbb{E} \left[\left\| \sum_{i=1}^n \varepsilon_i x_i \right\|_\infty \right].$$

Now, each coordinate j of $\sum_{i=1}^n \varepsilon_i x_i$ is $\sum_{i=1}^n x_{ij}^2$ -sub-Gaussian, and thus using that $\mathbb{E}[\max_{j \leq d} Z_j] \leq \sqrt{2\sigma^2 \log d}$ for arbitrary σ^2 -sub-Gaussian Z_j (see Exercise 4.9), we have

$$R_n(\mathcal{F} \mid x_1^n) \leq \frac{r}{n} \sqrt{2 \log(2d) \max_j \sum_{i=1}^n x_{ij}^2}.$$

To facilitate comparison with Example 5.2.8, suppose that the vectors x_i all satisfy $\|x_i\|_\infty \leq b$. In this case, the preceding inequality implies that $R_n(\mathcal{F} \mid x_1^n) \leq rb\sqrt{2 \log(2d)/\sqrt{n}}$. In contrast, the ℓ_2 -norm of such x_i may satisfy $\|x_i\|_2 = b\sqrt{d}$, so that the bounds of Example 5.2.7 scale instead as $rb\sqrt{d}/\sqrt{n}$, which can be exponentially larger. \diamond

These examples are sufficient to derive a few sophisticated risk bounds. We focus on the case where we have a loss function applied to some class with reasonable Rademacher complexity, in which case it is possible to recenter the loss class and achieve reasonable complexity bounds. The coming proposition does precisely this in the case of margin-based binary classification. Consider points $(x, y) \in \mathcal{X} \times \{\pm 1\}$, and let \mathcal{F} be an arbitrary class of functions $f : \mathcal{X} \rightarrow \mathbb{R}$ and $\mathcal{L} = \{(x, y) \mapsto \ell(yf(x))\}_{f \in \mathcal{F}}$ be the induced collection of losses. As a typical example, we might have $\ell(t) = [1 - t]_+$, $\ell(t) = e^{-t}$, or $\ell(t) = \log(1 + e^{-t})$. We have the following proposition.

Proposition 5.2.9. *Let \mathcal{F} and \mathcal{X} be such that $\sup_{x \in \mathcal{X}} |f(x)| \leq M$ for $f \in \mathcal{F}$ and assume that ℓ is L -Lipschitz. Define the empirical and population risks $\hat{L}_n(f) := P_n \ell(Yf(X))$ and $L(f) := P \ell(Yf(X))$. Then*

$$\mathbb{P} \left(\sup_{f \in \mathcal{F}} |\hat{L}_n(f) - L(f)| \geq 4LR_n(\mathcal{F}) + t \right) \leq 2 \exp \left(-\frac{nt^2}{2L^2M^2} \right) \quad \text{for } t \geq 0.$$

Proof We may recenter the class \mathcal{L} , that is, replace $\ell(\cdot)$ with $\ell(\cdot) - \ell(0)$, without changing $\hat{L}_n(f) - L(f)$. Call this class \mathcal{L}_0 , so that $\|P_n - P\|_{\mathcal{L}} = \|P_n - P\|_{\mathcal{L}_0}$. This recentered class satisfies bounded differences with constant $2ML$, as $|\ell(yf(x)) - \ell(y'f(x'))| \leq L|yf(x) - y'f(x')| \leq 2LM$, as in the proof of Proposition 5.1.1. Applying Proposition 5.1.1 and then Corollary 5.1.3 and gives

that $\mathbb{P}(\sup_{f \in \mathcal{F}} |\widehat{L}_n(f) - L(f)| \geq 2R_n(\mathcal{L}_0) + t) \leq \exp(-\frac{nt^2}{2M^2L^2})$ for $t \geq 0$. Then applying the contraction inequality (Theorem 5.1.6) yields $R_n(\mathcal{L}_0) \leq 2LR_n(\mathcal{F})$, giving the result. \square

Let us give a few example applications of these ideas.

Example 5.2.10 (Support vector machines and hinge losses): In the support vector machine problem, we receive data $(X_i, Y_i) \in \mathbb{R}^d \times \{\pm 1\}$, and we seek to minimize average of the losses $\ell(\theta; (x, y)) = [1 - y\theta^T x]_+$. We assume that the space \mathcal{X} has $\|x\|_2 \leq b$ for $x \in \mathcal{X}$ and that $\Theta = \{\theta \in \mathbb{R}^d \mid \|\theta\|_2 \leq r\}$. Applying Proposition 5.2.9 gives

$$\mathbb{P}\left(\sup_{\theta \in \Theta} |P_n \ell(\theta; (X, Y)) - P \ell(\theta; (X, Y))| \geq 4R_n(\mathcal{F}_\Theta) + t\right) \leq \exp\left(-\frac{nt^2}{2r^2b^2}\right),$$

where $\mathcal{F}_\Theta = \{f_\theta(x) = \theta^T x\}_{\theta \in \Theta}$. Now, we apply Example 5.2.7, which implies that

$$R_n(\phi \circ \mathcal{F}_\Theta) \leq 2R_n(\mathcal{F}_\Theta) \leq \frac{2rb}{\sqrt{n}}.$$

That is, we have

$$\mathbb{P}\left(\sup_{\theta \in \Theta} |P_n \ell(\theta; (X, Y)) - P \ell(\theta; (X, Y))| \geq \frac{4rb}{\sqrt{n}} + t\right) \leq \exp\left(-\frac{nt^2}{2(rb)^2}\right),$$

so that P_n and P become close at rate roughly rb/\sqrt{n} in this case. \diamond

Example 5.2.10 is what is sometimes called a “dimension free” convergence result—there is no explicit dependence on the dimension d of the problem, except as the radii r and b make explicit. One consequence of this is that if x and θ instead belong to a Hilbert space (potential infinite dimensional) with inner product $\langle \cdot, \cdot \rangle$ and norm $\|x\|^2 = \langle x, x \rangle$, but for which we are guaranteed that $\|\theta\| \leq r$ and similarly $\|x\| \leq b$, then the result still applies. Extending this to other function classes is reasonably straightforward, and we present a few examples in the exercises.

When we do not have the simplifying structure of $\ell(yf(x))$ identified in the preceding examples, we can still provide guarantees of generalization using the covering number guarantees introduced in Section 5.1.2. The most common and important case is when we have a Lipschitzian loss function in an underlying parameter θ .

Example 5.2.11 (Lipschitz functions over a norm-bounded parameter space): Consider the parametric loss minimization problem

$$\underset{\theta \in \Theta}{\text{minimize}} \quad L(\theta) := \mathbb{E}[\ell(\theta; Z)]$$

for a loss function ℓ that is M -Lipschitz (with respect to the norm $\|\cdot\|$) in its argument, where for normalization we assume $\inf_{\theta \in \Theta} \ell(\theta, z) = 0$ for each z . Then the metric entropy of Θ bounds the metric entropy of the loss class $\mathcal{F} := \{z \mapsto \ell(\theta, z)\}_{\theta \in \Theta}$ for the supremum norm $\|\cdot\|_\infty$. Indeed, for any pair θ, θ' , we have

$$\sup_z |\ell(\theta, z) - \ell(\theta', z)| \leq M \|\theta - \theta'\|,$$

and so an ϵ -cover of Θ is an $M\epsilon$ -cover of \mathcal{F} in supremum norm. In particular,

$$N(\epsilon, \mathcal{F}, \|\cdot\|_\infty) \leq N(\epsilon/M, \Theta, \|\cdot\|).$$

Assume that $\Theta \subset \{\theta \mid \|\theta\| \leq b\}$ for some finite b . Then Lemma 5.1.10 guarantees that $\log N(\epsilon, \Theta, \|\cdot\|) \leq d \log(1 + 2/\epsilon) \lesssim d \log \frac{1}{\epsilon}$, and so the classical covering number argument in Example 5.1.9 gives

$$\mathbb{P} \left(\sup_{\theta \in \Theta} |P_n \ell(\theta, Z) - P \ell(\theta, Z)| \geq t \right) \leq \exp \left(-c \frac{nt^2}{b^2 M^2} + Cd \log \frac{M}{t} \right),$$

where c, C are numerical constants. In particular, taking $t^2 \asymp \frac{M^2 b^2 d}{n} \log \frac{n}{\delta}$ gives that

$$|P_n \ell(\theta, Z) - P \ell(\theta, Z)| \lesssim \frac{Mb \sqrt{d \log \frac{n}{\delta}}}{\sqrt{n}}$$

with probability at least $1 - \delta$. \diamond

5.2.3 Structural risk minimization and adaptivity

In general, for a given function class \mathcal{F} , we can always decompose the excess risk into the *approximation/estimation* error decomposition. That is, let

$$L^* = \inf_f L(f),$$

where the preceding infimum is taken across *all* (measurable) functions. Then we have

$$L(\hat{f}_n) - L^* = \underbrace{L(\hat{f}_n) - \inf_{f \in \mathcal{F}} L(f)}_{\text{estimation}} + \underbrace{\inf_{f \in \mathcal{F}} L(f) - L^*}_{\text{approximation}}. \quad (5.2.6)$$

There is often a tradeoff between these two, analogous to the bias/variance tradeoff in classical statistics; if the approximation error is very small, then it is likely hard to guarantee that the estimation error converges quickly to zero, while certainly a constant function will have low estimation error, but may have substantial approximation error. With that in mind, we would like to develop procedures that, rather than simply attaining good performance for the class \mathcal{F} , are guaranteed to trade-off in an appropriate way between the two types of error. This leads us to the idea of *structural risk minimization*.

In this scenario, we assume we have a sequence of classes of functions, $\mathcal{F}_1, \mathcal{F}_2, \dots$, of increasing complexity, meaning that $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots$. For example, in a linear classification setting with vectors $x \in \mathbb{R}^d$, we might take a sequence of classes allowing increasing numbers of non-zeros in the classification vector θ :

$$\mathcal{F}_1 := \left\{ f_\theta(x) = \theta^\top x \text{ such that } \|\theta\|_0 \leq 1 \right\}, \mathcal{F}_2 := \left\{ f_\theta(x) = \theta^\top x \text{ such that } \|\theta\|_0 \leq 2 \right\}, \dots$$

More broadly, let $\{\mathcal{F}_k\}_{k \in \mathbb{N}}$ be a (possibly infinite) increasing sequence of function classes. We assume that for each \mathcal{F}_k and each $n \in \mathbb{N}$, there exists a constant $C_{n,k}(\delta)$ such that we have the uniform generalization guarantee

$$\mathbb{P} \left(\sup_{f \in \mathcal{F}_k} \left| \hat{L}_n(f) - L(f) \right| \geq C_{n,k}(\delta) \right) \leq \delta \cdot 2^{-k}.$$

For example, by Corollary 5.2.6, if \mathcal{F} is finite we may take

$$C_{n,k}(\delta) = \sqrt{2\sigma^2 \frac{\log |\mathcal{F}_k| + \log \frac{1}{\delta} + k \log 2}{n}}.$$

(We will see in subsequent sections of the course how to obtain other more general guarantees.)

We consider the following *structural risk minimization* procedure. First, given the empirical risk \hat{L}_n , we find the model collection \hat{k} minimizing the penalized risk

$$\hat{k} := \operatorname{argmin}_{k \in \mathbb{N}} \left\{ \inf_{f \in \mathcal{F}_k} \hat{L}_n(f) + C_{n,k}(\delta) \right\}. \quad (5.2.7a)$$

We then choose \hat{f} to minimize the risk over the estimated “best” class $\mathcal{F}_{\hat{k}}$, that is, set

$$\hat{f} := \operatorname{argmin}_{f \in \mathcal{F}_{\hat{k}}} \hat{L}_n(f). \quad (5.2.7b)$$

With this procedure, we have the following theorem.

Theorem 5.2.12. *Let \hat{f} be chosen according to the procedure (5.2.7a)–(5.2.7b). Then with probability at least $1 - \delta$, we have*

$$L(\hat{f}) \leq \inf_{k \in \mathbb{N}} \inf_{f \in \mathcal{F}_k} \{L(f) + 2C_{n,k}(\delta)\}.$$

Proof First, we have by the assumed guarantee on $C_{n,k}(\delta)$ that

$$\begin{aligned} & \mathbb{P} \left(\exists k \in \mathbb{N} \text{ and } f \in \mathcal{F}_k \text{ such that } \sup_{f \in \mathcal{F}_k} |\hat{L}_n(f) - L(f)| \geq C_{n,k}(\delta) \right) \\ & \leq \sum_{k=1}^{\infty} \mathbb{P} \left(\exists f \in \mathcal{F}_k \text{ such that } \sup_{f \in \mathcal{F}_k} |\hat{L}_n(f) - L(f)| \geq C_{n,k}(\delta) \right) \leq \sum_{k=1}^{\infty} \delta \cdot 2^{-k} = \delta. \end{aligned}$$

On the event that $\sup_{f \in \mathcal{F}_k} |\hat{L}_n(f) - L(f)| < C_{n,k}(\delta)$ for all k , which occurs with probability at least $1 - \delta$, we have

$$L(\hat{f}) \leq \hat{L}_n(\hat{f}) + C_{n,\hat{k}}(\delta) = \inf_{k \in \mathbb{N}} \inf_{f \in \mathcal{F}_k} \left\{ \hat{L}_n(f) + C_{n,k}(\delta) \right\} \leq \inf_{k \in \mathbb{N}} \inf_{f \in \mathcal{F}_k} \{L(f) + 2C_{n,k}(\delta)\}$$

by our choice of \hat{f} . This is the desired result. \square

We conclude with a final example, using our earlier floating point bound from Example 5.2.5, coupled with Corollary 5.2.6 and Theorem 5.2.12.

Example 5.2.13 (Structural risk minimization with floating point classifiers): Consider again our floating point example, and let the function class \mathcal{F}_k consist of functions defined by at most k double-precision floating point values, so that $\log |\mathcal{F}_k| \leq 45d$. Then by taking

$$C_{n,k}(\delta) = \sqrt{\frac{\log \frac{1}{\delta} + 65k \log 2}{2n}}$$

we have that $|\widehat{L}_n(f) - L(f)| \leq C_{n,k}(\delta)$ simultaneously for all $f \in \mathcal{F}_k$ and all \mathcal{F}_k , with probability at least $1 - \delta$. Then the empirical risk minimization procedure (5.2.7) guarantees that

$$L(\widehat{f}) \leq \inf_{k \in \mathbb{N}} \left\{ \inf_{f \in \mathcal{F}_k} L(f) + \sqrt{\frac{2 \log \frac{1}{\delta} + 91k}{n}} \right\}.$$

Roughly, we trade between small risk $L(f)$ —as the risk $\inf_{f \in \mathcal{F}_k} L(f)$ must be decreasing in k —and the estimation error penalty, which scales as $\sqrt{(k + \log \frac{1}{\delta})/n}$. \diamond

5.3 M-estimators and estimation

In many problems in statistics and machine learning, we seek not just to have small loss on some future data but to actually recover parameters of interest. For example, in a regression problem where we model

$$y = \langle x, \theta \rangle + \varepsilon,$$

we often care about the actual values of θ . Less prosaically, in the latter part of the book we will develop a number of fundamental limits and lower bounds; it is good to have algorithms and upper bounds to demonstrate their tightness!

To that end, we here develop representative finite-sample results on the convergence of different estimators. We focus on *M-estimators*, meaning those that arise from minimization of a loss function $\ell(\theta, z)$ convex in θ , where z is problem data. Exponential families, from Chapter 3, provide a natural examples here.

Example 5.3.1 (Exponential families and log loss): Let $p_\theta(x) = \exp(\langle \theta, \phi(x) \rangle - A(\theta))$ be the density of an exponential family. Then $\ell(\theta, x) := -\log p_\theta(x) = -\langle \theta, \phi(x) \rangle + A(\theta)$ is convex and \mathcal{C}^∞ over $\Theta := \{\theta \mid A(\theta) < \infty\}$. \diamond

Example 5.3.2 (Logistic regression): In binary logistic regression (Example 3.4.2) with labels $y \in \{\pm 1\}$, we have data $z = (x, y) \in \mathbb{R}^d \times \{\pm 1\}$, and the loss

$$\ell(\theta, z) = -\log p_\theta(y \mid x) = \log \left(1 + \exp(-yx^\top \theta) \right),$$

which is \mathcal{C}^∞ , domain all of \mathbb{R}^d , and Lipschitz-continuous derivatives of all orders. (Though the particular Lipschitz constant depends on x). \diamond

Regardless, we will consider general rates of convergence and estimation error for estimators minimizing the empirical loss

$$L_n(\theta) := P_n \ell(\theta, Z) = \frac{1}{n} \sum_{i=1}^n \ell(\theta, Z_i),$$

which approximates the population loss

$$L(\theta) := \mathbb{E}_P[\ell(\theta, Z)],$$

where $Z_i \stackrel{\text{iid}}{\sim} P$ and $P_n = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{Z_i}$ is the usual empirical distribution. Thus, for a closed convex set $\Theta \subset \mathbb{R}^d$, we will study the *M-estimator*

$$\widehat{\theta}_n := \operatorname{argmin}_{\theta \in \Theta} L_n(\theta), \tag{5.3.1}$$

providing prototypical arguments for its convergence. Often, we will take $\Theta = \mathbb{R}^d$, though this will not be essential.

Based on the results in the preceding sections on uniform convergence, one natural idea is to use uniform convergence: if we can argue that

$$\sup_{\theta \in \Theta} |L_n(\theta) - L(\theta)| \rightarrow 0,$$

then so long as the minimizer θ^* of L is unique and $L(\theta) - L(\theta^*)$ grows with the distance $\|\theta - \theta^*\|$, we necessarily have $\hat{\theta}_n \rightarrow \theta^*$. Unfortunately, this naive approach typically fails to achieve the correct convergence rates, let alone the correct dependence on problem parameters. We therefore take another approach.

JCD Comment: Will need some figures here / illustrations

Recall that a twice differentiable function L is convex if and only if $\nabla^2 L(\theta) \succeq 0$ for all $\theta \in \text{dom } L$. We thus expect that L should have some *quadratic growth* around its minimizer θ^* , meaning that in a neighborhood of θ^* , we have

$$L(\theta) \geq L(\theta^*) + \frac{\lambda}{2} \|\theta - \theta^*\|_2^2$$

for θ near enough θ^* . In such a situation, because the sampled L_n is also convex and approximates L , we then expect that for parameters θ far enough from θ^* that the growth of $L(\theta)$ above $L(\theta^*)$ dominates the noise inherent in the sampling, we necessarily have $L_n(\theta) > L_n(\theta^*)$. Because the empirical minimizer $\hat{\theta}_n$ necessarily satisfies $L_n(\hat{\theta}_n) \leq L_n(\theta^*)$, we would never choose such a distant parameter, thus implying a convergence rate. To make this type of argument rigorous requires a bit of convex analysis and sampling theory; luckily, we are by now well-equipped to address this.

5.3.1 Standard conditions and convex optimization

To provide relatively clean results, we will consider a collection of loss functions to simplify analysis. As we will see, the assumed conditions are not too onerous, as many families of losses satisfy them. We can relax them using some of the more sophisticated concentration inequalities we have developed. We therefore make the following standing assumption.

Assumption A.5.1 (Standard conditions). *For each $z \in \mathcal{Z}$, the losses $\ell(\theta, z)$ are convex in θ . There are constants $M_0, M_1, M_2 < \infty$ such that for each $z \in \mathcal{Z}$,*

(i) $\|\nabla \ell(\theta^*, z)\|_2 \leq M_0$

(ii) $\|\nabla^2 \ell(\theta^*, z)\|_{\text{op}} \leq M_1$, and

(iii) the Hessian $\nabla^2 \ell(\theta, z)$ is M_2 -Lipschitz continuous in a neighborhood of radius $r > 0$ around θ^* , meaning $\|\nabla^2 \ell(\theta_0, z) - \nabla^2 \ell(\theta_1, z)\|_{\text{op}} \leq M_0 \|\theta_0 - \theta_1\|_2$ whenever $\|\theta_i - \theta^*\|_2 \leq r$.

Additionally, the minimizer $\theta^* = \arg\min_{\theta} L(\theta)$ exists and for a $\lambda > 0$ satisfies

$$\nabla^2 L(\theta^*) \succeq \lambda I.$$

The “standard conditions” in Assumption A.5.1 are not so onerous. As we see when we specialize our coming results to exponential family models in Section 5.3.4, Assumption A.5.1 holds essentially as soon as the family is minimal and $\phi(x)$ is bounded. The existence of minimizers can be somewhat more subtle to guarantee than the smoothness conditions (i)–(iii), though these are typically straightforward. (For more on the existence of minimizers, see Exercise 5.10.)

To quickly highlight the conditions, we revisit binary logistic and robust regression.

Example (Example 5.3.2 continued): For logistic regression with labels $y \in \{\pm 1\}$, we have

$$\nabla \ell(\theta, (x, y)) = -\frac{1}{1 + e^{yx^\top \theta}} yx \quad \text{and} \quad \nabla^2 \ell(\theta, (x, y)) = p_\theta(y | x)(1 - p_\theta(y | x))xx^\top.$$

Then Assumptions (i)–(iii) hold so long as $\sup_{x \in \mathcal{X}} \|x\|_2 < \infty$, with $M_0 = \sup_{x \in \mathcal{X}} \|x\|_2$ and $M_1 = \frac{1}{4}M_0^2$. We revisit the existence of minimizers in the sequel, noting that because $0 < p_\theta(y | x)(1 - p_\theta(y | x)) \leq \frac{1}{4}$ for any θ, x, y , if a minimizer exists then it is necessarily unique as soon as $\mathbb{E}[XX^\top] \succ 0$. \diamond

Example 5.3.3 (Robust regression): In robust regression, we wish to best approximate responses $y \in \mathbb{R}$ via linear functions $x^\top \theta$, but because of outliers, do not use the squared loss. Thus, for a smooth symmetric convex function h with bounded derivatives, we take

$$\ell(\theta, (x, y)) = h(\langle x, \theta \rangle - y),$$

so that

$$\nabla \ell(\theta, (x, y)) = h'(\langle x, \theta \rangle - y)x \quad \text{and} \quad \nabla^2 \ell(\theta, (x, y)) = h''(\langle x, \theta \rangle - y)xx^\top.$$

A prototypical example is $h(t) = \log(1 + e^t) + \log(1 + e^{-t})$, which satisfies $h'(t) = \frac{e^t - 1}{e^t + 1} \in [-1, 1]$ and $h''(t) = \frac{2e^t}{(e^t + 1)^2} \in (0, 1/2]$. So long as the covariates $x \in \mathcal{X}$ have finite radius $\text{rad}_2(\mathcal{X}) := \sup_{x \in \mathcal{X}} \|x\|_2$, we obtain the Lipschitz constant bounds

$$M_0 \leq \text{rad}_2(\mathcal{X}), \quad M_1 \lesssim \text{rad}_2^2(\mathcal{X}), \quad \text{and} \quad M_2 \lesssim \text{rad}_2^3(\mathcal{X})$$

for Assumption A.5.1, parts (i)–(iii). In general, if h is symmetric with $h''(0) > 0$, then minimizers exist whenever Y is non-pathological and $\mathbb{E}[XX^\top] \succ 0$. Exercise 5.11 asks you to prove this last claim on existence of minimizers. \diamond

5.3.2 Some growth properties of convex functions

As we discuss above, we will roughly proceed in our analysis by showing that the growth of the loss function dominates the noise inherent in sampling. To do so, we will rely on certain growth properties of convex functions. We collect them here, as they provide the fundamental building block for convergence analysis.

First, we show that for any convex function h , if there exists a “shell” $S = \{\theta \mid \|\theta - \theta_0\|_2 = r\}$ around some point θ_0 for which $h(\theta) > h(\theta_0)$ for all $\theta \in S$, then necessarily the minimizer $\hat{\theta} = \text{argmin}_\theta h(\theta)$ satisfies $\|\hat{\theta} - \theta_0\|_2 < r$.

JCD Comment: Figure(s)

Lemma 5.3.4. *Let h be convex and $\theta_0 \in \text{dom } h$ and v an arbitrary vector. Then for all $t \geq 1$, $h(\theta_0 + tv) - h(\theta_0) \geq t(h(\theta_1) - h(\theta_0))$.*

Proof Let $\theta_t = \theta_0 + tv$ for $t \geq 1$. Then for $t \geq 1$, we can write $\theta_1 = \frac{1}{t}\theta_t + (1 - \frac{1}{t})\theta_0$, and so

$$h(\theta_1) = h\left(\left(1 - \frac{1}{t}\right)\theta_0 + \frac{1}{t}\theta_t\right) \leq \left(1 - \frac{1}{t}\right)h(\theta_0) + \frac{1}{t}h(\theta_t).$$

Rearranging yields

$$\left(1 - \frac{1}{t}\right)(h(\theta_1) - h(\theta_0)) \leq \frac{1}{t}(h(\theta_t) - h(\theta_1)),$$

and multiplying through by t and rearranging implies the desired result. \square

Extending this to shells, we have the following result.

Lemma 5.3.5. *Let h be convex and $\theta_0 \in \text{dom } h$ and assume that for some $\epsilon > 0$ and $\delta > 0$ that $h(\theta_0 + v) \geq h(\theta_0) + \delta$ for all $\|v\|_2 = \epsilon$. Then for all θ for which $\|\theta - \theta_0\|_2 \geq \epsilon$,*

$$h(\theta) - h(\theta_0) \geq \frac{\delta}{\epsilon} \|\theta - \theta_0\|_2.$$

Proof For any θ with $\|\theta - \theta_0\|_2 \geq \epsilon$, we can write $\theta = \theta_0 + tv$ for $v = \epsilon \frac{\theta - \theta_0}{\|\theta - \theta_0\|_2}$ and $t = \frac{\|\theta - \theta_0\|_2}{\epsilon}$. Apply Lemma 5.3.4 with the substitution $\delta \geq h(\theta_1) - h(\theta_0)$. \square

Now we connect these results to the growth of suitably smooth convex functions. Here, we wish to argue that the minimizer of a convex function h is not too far from a benchmark point θ_0 at which h has strong upward curvature and small gradient.

JCD Comment: Figure.

Lemma 5.3.6. *Let h be convex and $\lambda > 0$, $\gamma \geq 0$, and $\epsilon \geq \frac{2\gamma}{\lambda}$. Assume that for some θ_0 , we have $\|\nabla h(\theta_0)\|_2 \leq \gamma$ and $\nabla^2 h(\theta) \succeq \lambda I$ for all θ satisfying $\|\theta - \theta_0\|_2 \leq \epsilon$. Then the minimizer $\hat{\theta} = \text{argmin}_{\theta} h(\theta)$ exists and satisfies*

$$\|\hat{\theta} - \theta_0\|_2 \leq \frac{2\gamma}{\lambda}.$$

Proof By Taylor's theorem, for any θ we have

$$h(\theta) = h(\theta_0) + \langle \nabla h(\theta_0), \theta - \theta_0 \rangle + \frac{1}{2}(\theta - \theta_0)^\top \nabla^2 h(\bar{\theta})(\theta - \theta_0)$$

for a point $\bar{\theta}$ on the line between θ and θ_0 . Now, let us take θ such that $\|\theta - \theta_0\|_2 \leq t$. Then by assumption $\nabla^2 h(\bar{\theta}) \succeq \lambda I$, and so we have

$$\begin{aligned} h(\theta) &\geq h(\theta_0) + \langle \nabla h(\theta_0), \theta - \theta_0 \rangle + \frac{\lambda}{2} \|\theta - \theta_0\|_2^2 \\ &\geq h(\theta_0) - \gamma \|\theta - \theta_0\|_2 + \frac{\lambda}{2} \|\theta - \theta_0\|_2^2. \end{aligned}$$

by assumption that $\|\nabla h(\theta_0)\|_2 \leq \gamma$ and the Cauchy-Schwarz inequality.

Fix $t \geq 0$. If we can show that $h(\theta) > h(\theta_0)$ for all θ satisfying $\|\theta - \theta_0\|_2 = t$, then Lemma 5.3.5 implies that $h(\theta) > h(\theta_0)$ whenever $\|\theta - \theta_0\|_2 \geq t$, so that necessarily $\|\hat{\theta} - \theta_0\|_2 < t$. Returning to the previous display and letting $t = \|\theta - \theta_0\|_2$, note that

$$h(\theta) \geq h(\theta_0) - \gamma t + \frac{\lambda}{2} t^2,$$

and as $-\gamma t + \frac{\lambda}{2}t^2 = t(\frac{\lambda}{2}t - \gamma) > 0$ whenever $t > \frac{2\gamma}{\lambda}$. As by assumption $\nabla^2 h(\theta) \succeq \lambda I$ whenever $\|\theta - \theta_0\|_2 \leq \epsilon$ for some $\epsilon \geq \frac{2\gamma}{\lambda}$, this implies the result. \square

Lemma 5.3.7. *Let h be convex and assume that $\nabla^2 h$ is M_2 -Lipschitz (part (iii) of Assumption A.5.1). Let $\lambda > 0$ be large enough and $\gamma > 0$ be small enough that $\gamma < \frac{\lambda^2}{8M_2}$. Then if both $\nabla^2 h(\theta_0) \succeq \lambda I$ and $\|\nabla h(\theta_0)\|_2 \leq \gamma$, the minimizer $\hat{\theta} = \operatorname{argmin}_{\theta} h(\theta)$ exists and satisfies*

$$\|\hat{\theta} - \theta_0\|_2 \leq \frac{4\gamma}{\lambda}.$$

Proof By Lemma 5.3.6, it is enough to show that $\nabla^2 h(\theta) \succeq \frac{\lambda}{2}I$ for all θ with $\|\theta - \theta_0\|_2 \leq \frac{4\gamma}{\lambda}$. For this, we use the M_2 -Lipschitz continuity of $\nabla^2 h$ to obtain that for any θ with $\|\theta - \theta_0\|_2 = t$,

$$\nabla^2 h(\theta) \succeq \nabla^2 h(\theta_0) - M_2 \|\theta - \theta_0\|_2 I \succeq (\lambda - M_2 t)I.$$

So if $t \leq \frac{\lambda}{2M_2}$ we have $\nabla^2 h(\theta) \succeq \lambda I$. Because $\frac{4\gamma}{\lambda} \leq \frac{\lambda}{2M_2}$ by assumption, we have $\nabla^2 h(\theta) \succeq \frac{\lambda}{2}I$ whenever $\|\theta - \theta_0\|_2 \leq \frac{4\gamma}{\lambda}$, yielding the result. \square

5.3.3 Convergence analysis for convex M-estimators

By leveraging Lemma 5.3.7, to show a convergence rate guarantee for the empirical minimizer $\hat{\theta}_n$, it is evidently sufficient to demonstrate two (related) conditions: that for some sequence $\gamma_n \rightarrow 0$, we have

$$\|\nabla L_n(\theta^*)\|_2 \leq \gamma_n \tag{5.3.2a}$$

with high probability, and that for some $\lambda > 0$, we have

$$\nabla^2 L(\theta^*) \succeq \lambda I \tag{5.3.2b}$$

with high probability. Happily, the convergence guarantees we develop in Chapter 4 provide precisely the tools to do this.

Theorem 5.3.8. *Let Assumption A.5.1 hold for the M-estimation problem (5.3.1). Let $\delta \in (0, \frac{1}{2})$, and define*

$$\gamma_n(\delta) := \frac{M_0}{\sqrt{n}} \left(1 + \sqrt{\log \frac{1}{\delta}} \right) \quad \text{and} \quad \epsilon_n(\delta) := \max \left\{ \frac{2 \|\nabla^2 L(\theta^*)\|_{\text{op}}}{\sqrt{n}} \sqrt{\log \frac{2d}{\delta}}, \frac{4M_1}{3n} \log \frac{2d}{\delta} \right\}.$$

Then we have both $\|\nabla L_n(\theta^)\|_2 \leq \gamma_n(\delta)$ and $\|\nabla^2 L_n(\theta^*) - \nabla^2 L(\theta^*)\|_{\text{op}} \leq \epsilon_n(\delta)$ with probability at least $1 - 2\delta$. So long as $\epsilon_n(\delta) \leq \frac{\lambda}{4}$ and $\gamma_n(\delta) \leq \frac{9\lambda^2}{128M_2}$, then with the same probability,*

$$\|\hat{\theta}_n - \theta^*\|_2 \leq \frac{16}{3} \frac{\gamma_n(\delta)}{\lambda}.$$

We defer the proof of the theorem to Section 5.3.5, instead providing commentary and a few examples of its application. Ignoring the numerical constants, the theorem roughly states the following: once n is large enough that

$$n \gtrsim \frac{M_2^2 M_0^2}{\lambda^4} \log \frac{1}{\delta} \quad \text{and} \quad n \gtrsim \frac{\|\nabla^2 L(\theta^*)\|_{\text{op}}^2}{\lambda^2} \log \frac{d}{\delta}, \quad (5.3.3)$$

with probability at least $1 - \delta$ we have

$$\|\hat{\theta}_n - \theta^*\|_2 \lesssim \frac{M_0}{\lambda\sqrt{n}} \sqrt{\log \frac{1}{\delta}}. \quad (5.3.4)$$

These finite sample results are, at least for large n , order optimal, as we will develop in the coming sections on fundamental limits. Nonetheless, the conditions (5.3.3) are stronger than necessary, typically requiring that n be quite large. In the exercises, we explore a class of *quasi-self-concordant* losses, where the second derivative controls the third derivative, allowing more direct application of Lemma 5.3.6, which allows reducing this sample size requirement. (See Exercises 5.5 and 5.6).

Example 5.3.9 (Logistic regression, Example 5.3.2 continued): Recalling the logistic loss $\ell(\theta, (x, y)) = \log(1 + e^{-y\langle x, \theta \rangle})$ for $y \in \{\pm 1\}$ and $x \in \mathbb{R}^d$, assume the domain \mathcal{X} consists of vectors x with $\|x\|_2 \leq \sqrt{d}$. For example, if $\mathcal{X} \subset [-1, 1]^d$, this holds. In this case, $M_0 = \sqrt{d}$, while $M_1 \leq \frac{1}{4}d$ and $M_2 \lesssim d^{3/2}$. Assuming that the population Hessian $\nabla^2 L(\theta^*) = \mathbb{E}[p_{\theta^*}(Y | X)(1 - p_{\theta^*}(Y | X))XX^\top]$ has minimal eigenvalue $\lambda_{\min}(\nabla^2 L(\theta^*)) \gtrsim 1$, then the conclusions of Theorem 5.3.8 apply as soon as $n \gtrsim d^4 \log \frac{1}{\delta}$. \diamond

When n is large enough, the guarantee (5.3.4) allows us to also make the heuristic asymptotic expansions for the exponential family models in Section 3.2.1 hold in finite samples. Let the conclusions of Theorem 5.3.8 hold, so that $\|\nabla^2 L_n(\theta^*) - \nabla^2 L(\theta^*)\|_{\text{op}} \leq \epsilon_n(\delta)$ and so on. Then once we know that $\hat{\theta}_n$ exists, by a Taylor expansion we can write

$$\begin{aligned} 0 &= \nabla L_n(\hat{\theta}_n) = \nabla L_n(\theta^*) + (\nabla^2 L_n(\theta^*) + E_n)(\hat{\theta}_n - \theta^*) \\ &= \nabla L_n(\theta^*) + (\nabla^2 L(\theta^*) + E'_n)(\hat{\theta}_n - \theta^*), \end{aligned}$$

where E_n is an error matrix satisfying $\|E_n\|_{\text{op}} \leq M_2 \|\hat{\theta}_n - \theta^*\|_2$ and $E'_n = E_n + \nabla^2 L_n(\theta^*) - \nabla^2 L(\theta^*)$ satisfies $\|E'_n\|_{\text{op}} \leq \|E_n\|_{\text{op}} + \epsilon_n(\delta)$ by the triangle inequality and Theorem 5.3.8. Using the infinite series expansion of the inverse

$$(A + E)^{-1} = A^{-1} + \sum_{i=1}^{\infty} (-1)^i (A^{-1} E)^i A^{-1},$$

valid for $A \succ 0$ whenever $\|E\|_{\text{op}} < \lambda_{\min}(A)$ (see Exercise 5.4), we therefore have

$$\hat{\theta}_n - \theta^* = -(\nabla^2 L(\theta^*) + E'_n)^{-1} \nabla L_n(\theta^*) = -\nabla^2 L(\theta^*)^{-1} \nabla L_n(\theta^*) + R_n,$$

where the remainder vector R_n satisfies

$$\begin{aligned} \|R_n\|_2 &\lesssim \|\nabla^2 L(\theta^*)^{-1} E'_n \nabla^2 L(\theta^*)^{-1}\|_{\text{op}} \|\nabla L_n(\theta^*)\|_2 \\ &\lesssim \frac{1}{\lambda^2} \left(\frac{M_2 M_0 \sqrt{\log \frac{1}{\delta}}}{\lambda \sqrt{n}} + \epsilon_n(\delta) \right) \|\nabla L_n(\theta^*)\|_2 \lesssim \frac{\log \frac{d}{\delta}}{\lambda^2 n} \cdot \left(\frac{M_2 M_0^2}{\lambda} + \|\nabla^2 L(\theta^*)\|_{\text{op}} \right) \end{aligned}$$

with probability at least $1 - \delta$. We summarize this in the following corollary.

Corollary 5.3.10. *Let the conditions of Theorem 5.3.8 hold. Then there exists a problem-dependent constant C such that the following holds: for any $\delta > 0$, for any $n \geq C \log \frac{1}{\delta}$, with probability at least $1 - \delta$*

$$\hat{\theta}_n - \theta^* = -\nabla^2 L(\theta^*)^{-1} \nabla L_n(\theta^*) + R_n,$$

where the remainder R_n satisfies $\|R_n\|_2 \leq C \cdot \frac{1}{n} \log \frac{d}{\delta}$. The constant C may be taken to be continuous in all the problem parameters of Assumption A.5.1.

Corollary 5.3.10 highlights the two salient terms governing error in estimation problems: the curvature of the loss near the optimum, as $\nabla^2 L(\theta^*)$ contributes, and the variance in the gradients $\nabla L_n(\theta^*) = \frac{1}{n} \sum_{i=1}^n \nabla \ell(\theta^*, Z_i)$. When the Hessian term $\nabla^2 L(\theta^*)$ is “large,” meaning that $\nabla^2 L(\theta^*) \succeq \lambda I$ for some large value $\lambda > 0$, then estimation is easier: the curvature of the loss helps to identify θ^* . Conversely, when the variance $\text{Var}(\nabla \ell(\theta^*, Z)) = \mathbb{E}[\|\nabla \ell(\theta^*, Z)\|_2^2]$ is large, then estimation is more challenging. As a final remark, let us imagine that the remainder term R_n in the corollary, in addition to being small with high probability, satisfies $\mathbb{E}[\|R_n\|_2^2] \leq \frac{C}{n^2}$, where C is a problem-dependent constant. Let $G_n = -\nabla^2 L(\theta^*)^{-1} \nabla L_n(\theta^*)$ be the leading term in the expansion, which satisfies

$$\mathbb{E}[\|G_n\|_2^2] = \frac{1}{n} \text{Var}(L(\theta^*)^{-1} \nabla \ell(\theta^*, Z)) = \frac{\text{tr}(\nabla^2 L(\theta^*)^{-1} \text{Cov}(\nabla \ell(\theta^*, Z)) \nabla^2 L(\theta^*)^{-1})}{n}.$$

Then because $\|G_n + R_n\|_2^2 \leq \|G_n\|_2^2 + \|R_n\|_2^2 + 2\|G_n\|_2 \|R_n\|_2$, we have the heuristic

$$\begin{aligned} \mathbb{E}[\|\hat{\theta}_n - \theta^*\|_2^2] &= \mathbb{E}[\|G_n + R_n\|_2^2] \\ &= \mathbb{E}[\|G_n\|_2^2] + \mathbb{E}[\|R_n\|_2^2] \pm 2\mathbb{E}[\|G_n\|_2 \|R_n\|_2] \\ &\stackrel{(\star)}{=} \frac{\text{tr}(\nabla^2 L(\theta^*)^{-1} \text{Cov}(\nabla \ell(\theta^*, Z)) \nabla^2 L(\theta^*)^{-1})}{n} \pm \frac{C}{n^{3/2}}, \end{aligned} \quad (5.3.5)$$

where C is a problem-dependent constant and the step (\star) is heuristic. See Exercise 5.7 for one approach to make this step rigorous.

5.3.4 Consequences for exponential families and generalized linear models

Working through a few example applications of Corollary 5.3.10 with the exponential family and generalized linear models of Chapter 3 can help to make the results and connections clearer. Recall that for an exponential family model with loss $\ell(\theta; x) = -\log p_\theta(x) = -\langle \theta, \phi(x) \rangle + A(\theta)$, we heuristically derived in expression (3.2.4) that if the data were i.i.d. from the exponential family model P_{θ^*} , then

$$\hat{\theta}_n - \theta^* \sim \mathbf{N}(0, n^{-1} \cdot \nabla^2 A(\theta^*)^{-1}).$$

Corollary 5.3.10 presents one approach to make this rigorous. Assuming the sufficient statistics ϕ are bounded, we have $\nabla^2 L(\theta^*) = \nabla^2 A(\theta^*)$, and $\text{Cov}_{\theta^*}(\nabla L_n(\theta^*)) = \frac{1}{n} \text{Cov}_{\theta^*}(\phi(X)) = \frac{1}{n} \nabla^2 A(\theta^*)$. So

$$\hat{\theta}_n - \theta^* = -\nabla^2 A(\theta^*)^{-1} \frac{1}{n} \sum_{i=1}^n (\phi(X_i) - \nabla A(\theta^*)) + R_n,$$

where the remainder satisfies $\|R_n\|_2 \leq C \frac{1}{n} \log \frac{d}{\delta}$.

To obtain finite sample expected bounds requires a bit of tedium because of small probability events (e.g., that the sampled Hessian matrix $\nabla^2 L_n(\theta^*)$ fails to be invertible). One simple device

is to consider the estimator $\hat{\theta}_n$ only on some “good” event \mathcal{E}_n that occurs with high probability, for example, that the remainder R_n is small. The next corollary provides a prototypical result under the assumption that $X_i \stackrel{\text{iid}}{\sim} P_{\theta^*}$ for an exponential family model with bounded data $\sup_{x \in \mathcal{X}} \|\phi(x)\|_2 < \infty$ and positive definite Hessian $\nabla^2 A(\theta^*) \succ 0$.

Corollary 5.3.11. *Under the preceding conditions on the exponential family model P_{θ^*} , there exists a problem dependent constant $C < \infty$ such that the following holds: for any $k \geq 1$, there are events \mathcal{E}_n with probability $\mathbb{P}(\mathcal{E}_n) \geq 1 - \frac{1}{n^k}$ and*

$$\mathbb{E}_{\theta^*} \left[\|\hat{\theta}_n - \theta^*\|_2^2 \cdot \mathbf{1}\{\mathcal{E}_n\} \right] \leq \frac{1}{n} \text{tr}(\nabla^2 A(\theta^*)^{-1}) + \frac{Ck \log n}{n^{3/2}}.$$

The constant C may be taken continuous in θ^* .

Recalling the equality (3.3.2), we see that the Fisher information $\nabla^2 A(\theta)$ appears in a fundamental way for the exponential families. Proposition 3.5.1 shows that this quantity is fundamental, at least for testing; here it provides an upper bound on the convergence of the maximum likelihood estimator. Exercise 5.8 extends Corollary 5.3.11 to an equality to within lower order terms.

Proof Let $\delta = \delta_n > 0$ to be chosen and define the event \mathcal{E}_n to be that $\|R_n\|_2 \leq C \frac{1}{n} \log \frac{d}{\delta}$, which occurs with probability at least $1 - \delta$. Let

$$G_n = -\nabla^2 L(\theta^*)^{-1} \nabla L_n(\theta^*) = -\nabla^2 A(\theta^*)^{-1} P_n(\phi(X) - \nabla A(\theta^*))$$

be the mean-zero gradient term. Then $\hat{\theta}_n = G_n + R_n$, and $\|G_n + R_n\|_2^2 \leq \|G_n\|_2^2 + \|R_n\|_2^2 + 2\|G_n\|_2 \|R_n\|_2$. Then on the event \mathcal{E}_n we have $\|R_n\|_2^2 \leq \frac{C}{n} \log \frac{d}{\delta}$, and so

$$\mathbb{E} \left[\|\hat{\theta}_n - \theta^*\|_2^2 \mathbf{1}\{\mathcal{E}_n\} \right] \leq \mathbb{E}[\|G_n\|_2^2] + \mathbb{E}[\|R_n\|_2^2 \mathbf{1}\{\mathcal{E}_n\}] \mathbb{E}[\|G_n\|_2] + \frac{C^2}{n^2} \log^2 \frac{d}{\delta}.$$

Now note that $\mathbb{E}[\|G_n\|_2^2] = \frac{1}{n} \text{tr}(\nabla^2 A(\theta^*)^{-1})$, and set $\delta = \frac{1}{n^k}$. □

These ideas also extend to generalized linear models, such as linear, logistic, or Poisson regression (recall Chapter 3.4). For the abstract generalized linear model of predicting a target $y \in \mathcal{Y}$ from covariates $x \in \mathcal{X}$, we have

$$\ell(\theta, (x, y)) = -\log p_\theta(y | x) = -\phi(x, y)^\top \theta + A(\theta | x).$$

Because the log partition is \mathcal{C}^∞ , the smoothness conditions in Assumption A.5.1 then reduce to the boundedness

$$\text{rad}_2(\{\phi(x, y) | x \in \mathcal{X}, y \in \mathcal{Y}\}) := \sup_{x \in \mathcal{X}, y \in \mathcal{Y}} \|\phi(x, y)\|_2 < \infty.$$

Assuming that a minimizer $\theta^* = \text{argmin}_\theta L(\theta)$ exists, the (local) strong convexity condition that $\nabla^2 L(\theta^*) \succ 0$ then becomes that $\mathbb{E}[\nabla^2 A(\theta | X)] = \mathbb{E}[\text{Cov}_\theta(\phi(X, Y) | X)] \succ 0$. Exercise 5.10, part (c) gives general sufficient conditions for the existence of minimizers in GLMs.

For logistic regression (Example 3.4.2), these conditions correspond to a bound on the covariate data x , that $\mathbb{E}[XX^\top] \succ 0$, and that for each X , the label Y is non-deterministic. For Poisson regression (Example 3.4.4), we have $\ell(\theta, (x, y)) = -yx^\top \theta + e^{\theta^\top x}$. When the count data $Y \in \mathbb{N}$ can be unbounded, Assumption A.5.1.(i) may fail, because yx may be unbounded. If we *model* a data-generating process for which X and Y are both bounded using Poisson regression, however, then

the smoothness conditions in Assumption A.5.1 hold. (Again, see Exercise 5.10 for the existence of solutions.)

Regardless, by an argument completely parallel to that for Corollary 5.3.10, we can provide convergence rates for generalized linear model estimators. Here, we avoid the assumption of model fidelity, instead assuming that $\theta^* = \operatorname{argmin}_{\theta} L(\theta)$ exists and $\nabla^2 L(\theta^*) = \mathbb{E}[\nabla^2 A(\theta^* | X)] \succ 0$, so that θ^* is unique.

Corollary 5.3.12. *Let the preceding conditions hold and p_{θ} be a generalized linear model. Then there exists a problem constant $C < \infty$ such the following holds: for any $\delta \in (0, 1)$ and for all $n \geq C \log \frac{1}{\delta}$, with probability at least $1 - \delta$*

$$\hat{\theta}_n - \theta^* = -\mathbb{E}[\nabla^2 A(\theta^* | X)]^{-1} P_n(\phi(X, Y) - \nabla A(\theta^* | X)) + R_n,$$

where the remainder $\|R_n\|_2 \leq \frac{C}{n} \log \frac{1}{\delta}$.

When the generalized linear model P_{θ} is correct, so that $X \sim P$ marginally and $Y | X \sim P_{\theta}(\cdot | X)$, then $\operatorname{Cov}(\phi(X, Y) | X) = \nabla^2 A(\theta^* | X)$, and so in this case (again, for a sequence of events \mathcal{E}_n with probability at least $1 - 1/n^k$), we have

$$\mathbb{E}_{\theta^*} \left[\|\hat{\theta}_n - \theta^*\|_2^2 \mathbf{1}_{\{\mathcal{E}_n\}} \right] \leq \frac{\operatorname{tr}(\mathbb{E}[\nabla^2 A(\theta^* | X)]^{-1})}{n} + \frac{Ck \log n}{n^{3/2}}.$$

Note that this quantity is the trace of the inverse Fisher information (3.3.2) in the generalized linear model: the “larger” the information, the better estimation accuracy we can guarantee.

5.3.5 Proof of Theorem 5.3.8

The two key steps in the proof of the theorem are lemmas providing the guarantees (5.3.2).

Lemma 5.3.13. *Let Assumption A.5.1 hold. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$*

$$\|\nabla L_n(\theta^*)\|_2 \leq \frac{M_0}{\sqrt{n}} \left(1 + \sqrt{2 \log \frac{1}{\delta}} \right).$$

Proof The function $z_1^n \mapsto \|L_n(\theta)\|_2$ satisfies bounded differences: for any two empirical samples P_n, P'_n differing in only observation i ,

$$\begin{aligned} \left| \|P_n \nabla \ell(\theta, Z)\|_2 - \|P'_n \nabla \ell(\theta, Z)\|_2 \right| &\leq \|P_n \nabla \ell(\theta, Z) - P'_n \nabla \ell(\theta, Z)\|_2 \\ &\leq \frac{1}{n} \|\nabla \ell(\theta, Z_i) - \nabla \ell(\theta, Z'_i)\|_2 \leq \frac{2M_0}{n} \end{aligned}$$

by Assumption A.5.1.(i). Because $\mathbb{E}[\|\nabla L_n(\theta^*)\|_2] \leq \sqrt{\mathbb{E}[\|\nabla L_n(\theta^*)\|_2^2]} \leq M_0/\sqrt{n}$, Proposition 4.2.5 gives that

$$\mathbb{P}(\|\nabla L_n(\theta^*)\|_2 \geq M_0/\sqrt{n} + t) \leq \exp\left(-\frac{nt^2}{2M_0^2}\right)$$

for all $t \geq 0$. Solving for t in $\exp(-\frac{nt^2}{2M_0^2}) = \delta$ yields the lemma. \square

Lemma 5.3.14. *Let Assumption A.5.1 hold. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$*

$$\|\nabla^2 L_n(\theta^*) - \nabla^2 L(\theta^*)\|_{\text{op}} \leq \max \left\{ \frac{2 \|\nabla^2 L(\theta^*)\|_{\text{op}}}{\sqrt{n}} \sqrt{\log \frac{2d}{\delta}}, \frac{4M_1}{3n} \log \frac{2d}{\delta} \right\}.$$

Proof Because $\nabla^2 L_n(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla^2 \ell(\theta, Z_i)$ and $\|\nabla^2 \ell(\theta, z)\|_{\text{op}} \leq M_1$ by Assumption A.5.1.ii, Theorem 4.3.6 implies that

$$\mathbb{P} \left(\|\nabla^2 L_n(\theta^*) - \nabla^2 L(\theta^*)\|_{\text{op}} \geq t \right) \leq 2d \exp \left(- \min \left\{ \frac{nt^2}{4 \|\nabla^2 L(\theta^*)\|_{\text{op}}^2}, \frac{3t}{4M_1} \right\} \right).$$

Setting $t = \max \left\{ \frac{2 \|\nabla^2 L(\theta^*)\|_{\text{op}}}{\sqrt{n}} \sqrt{\log \frac{2d}{\delta}}, \frac{4M_1}{3n} \log \frac{2d}{\delta} \right\}$ gives the lemma. \square

Let $\epsilon_n(\delta)$ be the bound on the right side of Lemma 5.3.14. Then with probability at least $1 - \delta$,

$$\nabla^2 L_n(\theta^*) \succeq L(\theta^*) - \epsilon_n(\delta) I_d \geq \frac{3\lambda}{4} I_d$$

by the assumption that n is large enough and $\nabla^2 L(\theta^*) \succeq \lambda I$; we therefore have the first condition of Lemma 5.3.7 where $3\lambda/4$ replaces λ . Lemma 5.3.13 gives that with probability at least $1 - \delta$, $\|\nabla L_n(\theta^*)\|_2 \leq \gamma_n(\delta)$. Now use the assumption that $\gamma_n(\delta) \leq \frac{9\lambda^2}{128M_2}$, so that Lemma 5.3.7 implies

$$\|\hat{\theta}_n - \theta^*\|_2 \leq \frac{16}{3} \frac{\gamma_n(\delta)}{\lambda},$$

which proves the theorem.

5.4 Exercises

JCD Comment: Exercise ideas around this: We could try to do things with moment bounds. Like we'd use something like the Marcinkiewicz bounds, and then some moment bounds on matrices from different exercises, and we could say something.

Probably also write some moment bound guarantees for matrices with operator norms would be really neat.

Also, exercises that handle dimension scaling could be fun, along with (associated) convergence rates.

Exercise 5.1: In this question, we show how to use Bernstein-type (sub-exponential) inequalities to give sharp convergence guarantees. Recall (Example 4.1.14, Corollary 4.1.18, and inequality (4.1.6)) that if X_i are independent bounded random variables with $|X_i - \mathbb{E}[X]| \leq b$ for all i and $\text{Var}(X_i) \leq \sigma^2$, then

$$\max \left\{ \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n X_i \geq \mathbb{E}[X] + t \right), \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n X_i \leq \mathbb{E}[X] - t \right) \right\} \leq \exp \left(-\frac{1}{2} \min \left\{ \frac{5nt^2}{6\sigma^2}, \frac{nt}{2b} \right\} \right).$$

We consider minimization of loss functions ℓ over finite function classes \mathcal{F} with $\ell \in [0, 1]$, so that if $L(f) = \mathbb{E}[\ell(f, Z)]$ then $|\ell(f, Z) - L(f)| \leq 1$. Throughout this question, we let

$$L^* = \min_{f \in \mathcal{F}} L(f) \quad \text{and} \quad f^* \in \arg\min_{f \in \mathcal{F}} L(f).$$

We will show that, roughly, a procedure based on picking an empirical risk minimizer is unlikely to choose a function $f \in \mathcal{F}$ with bad performance, so that we obtain faster concentration guarantees.

(a) Argue that for any $f \in \mathcal{F}$

$$\mathbb{P}\left(\widehat{L}(f) \geq L(f) + t\right) \vee \mathbb{P}\left(\widehat{L}(f) \leq L(f) - t\right) \leq \exp\left(-\frac{1}{2} \min\left\{\frac{5}{6} \frac{nt^2}{L(f)(1-L(f))}, \frac{nt}{2}\right\}\right).$$

(b) Define the set of “bad” prediction functions $\mathcal{F}_{\epsilon \text{ bad}} := \{f \in \mathcal{F} : L(f) \geq L^* + \epsilon\}$. Show that for any fixed $\epsilon \geq 0$ and any $f \in \mathcal{F}_{2\epsilon \text{ bad}}$, we have

$$\mathbb{P}\left(\widehat{L}(f) \leq L^* + \epsilon\right) \leq \exp\left(-\frac{1}{2} \min\left\{\frac{5}{6} \frac{n\epsilon^2}{L^*(1-L^*) + \epsilon(1-\epsilon)}, \frac{n\epsilon}{2}\right\}\right).$$

(c) Let $\widehat{f}_n \in \operatorname{argmin}_{f \in \mathcal{F}} \widehat{L}(f)$ denote the empirical minimizer over the class \mathcal{F} . Argue that it is likely to have good performance, that is, for all $\epsilon \geq 0$ we have

$$\mathbb{P}\left(L(\widehat{f}_n) \geq L(f^*) + 2\epsilon\right) \leq \operatorname{card}(\mathcal{F}) \cdot \exp\left(-\frac{1}{2} \min\left\{\frac{5}{6} \frac{n\epsilon^2}{L^*(1-L^*) + \epsilon(1-\epsilon)}, \frac{n\epsilon}{2}\right\}\right).$$

(d) Using the result of part (c), argue that with probability at least $1 - \delta$,

$$L(\widehat{f}_n) \leq L(f^*) + \frac{4 \log \frac{|\mathcal{F}|}{\delta}}{n} + \sqrt{\frac{12}{5}} \cdot \frac{\sqrt{L^*(1-L^*) \cdot \log \frac{|\mathcal{F}|}{\delta}}}{\sqrt{n}}.$$

Why is this better than an inequality based purely on the boundedness of the loss ℓ , such as Theorem 5.2.4 or Corollary 5.2.6? What happens when there is a perfect risk minimizer f^* ?

Exercise 5.2: Prove Lemma 5.1.8.

Exercise 5.3: Consider a binary classification problem with logistic loss $\ell(\theta; (x, y)) = \log(1 + \exp(-y\theta^T x))$, where $\theta \in \Theta := \{\theta \in \mathbb{R}^d \mid \|\theta\|_1 \leq r\}$ and $y \in \{\pm 1\}$. Assume additionally that the space $\mathcal{X} \subset \{x \in \mathbb{R}^d \mid \|x\|_\infty \leq b\}$. Define the empirical and population risks $\widehat{L}_n(\theta) := P_n \ell(\theta; (X, Y))$ and $L(\theta) := P \ell(\theta; (X, Y))$, and let $\widehat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} \widehat{L}_n(\theta)$. Show that with probability at least $1 - \delta$ over $(X_i, Y_i) \stackrel{\text{iid}}{\sim} P$,

$$L(\widehat{\theta}_n) \leq \inf_{\theta \in \Theta} L(\theta) + C \frac{rb \sqrt{\log \frac{d}{\delta}}}{\sqrt{n}}$$

where $C < \infty$ is a numerical constant (you need not specify this).

Exercise 5.4: Let $A \succ 0$ be a positive definite matrix, and let E be Hermitian and satisfy $\|E\|_{\text{op}} < \lambda_{\min}(A)$. Define $S_k := A^{-1} + \sum_{i=1}^k (-1)^i (A^{-1}E)^i A^{-1}$.

(a) Show that for any $k \in \mathbb{N}$,

$$(A + E)S_k = I + (-1)^k (EA^{-1})^{k+1}.$$

(b) Argue that $S_\infty := \lim_k S_k$ exists and $(A + E)^{-1} = S_\infty = A^{-1} + \sum_{i=1}^{\infty} (-1)^i (A^{-1}E)^i A^{-1}$.

(c) Let $\gamma = \|E\|_{\text{op}}/\lambda_{\min}(A) < 1$. Show that

$$\|(A + E)^{-1} - S_k\|_{\text{op}} \leq \gamma^{k+1} \frac{1}{\lambda_{\min}(A)(1 - \gamma)}.$$

Exercise 5.5: Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a three-times differentiable convex function. We say f is *C-quasi-self-concordant* (q.s.c.) if $|f'''(t)| \leq C|f''(t)|$ for all $t \in \mathbb{R}$.

(a) Define $g(t) = \log f''(t)$. Show that if f is C-q.s.c., then $|g'(t)| \leq C$ and so for any $s \in \mathbb{R}$,

$$e^{-C|s|} f''(t) \leq f''(t + s) \leq e^{C|s|} f''(t).$$

(b) Show that the function $f(t) = \log(1 + e^t) + \log(1 + e^{-t})$ is q.s.c., and give its self-concordance parameter.

(c) Show that the function $f(t) = \log(1 + e^t)$ is q.s.c., and give its self-concordance parameter.

(d) Let f be C-q.s.c. and for a fixed $x \in \mathbb{R}^d$, define $h(\theta) := f(\langle \theta, x \rangle)$. Show that for any θ, θ_0 with $\Delta = \theta - \theta_0$, h satisfies

$$e^{-C|\langle \Delta, x \rangle|} \nabla^2 h(\theta) \preceq \nabla^2 h(\theta_0) \preceq e^{C|\langle \Delta, x \rangle|} \nabla^2 h(\theta).$$

Exercise 5.6 (Quasi self-concordant M-estimators [152]): Consider a prediction problem of predicting targets y from vectors $x \in \mathbb{R}^d$. A loss ℓ is a *C-quasi self-concordant* loss if we can write

$$\ell(\theta, (x, y)) = h(\langle \theta, x \rangle, y),$$

where for each y , $h(\cdot, y)$ is a C-q.s.c. function (recall Exercise 5.5).

(a) Show that logistic regression with loss $\ell(\theta, (x, y)) = \log(1 + e^{-y\langle x, \theta \rangle})$ and robust linear regression with loss $\ell(\theta, (x, y)) = \log(1 + e^{y\langle x, \theta \rangle}) + \log(1 + e^{\langle x, \theta \rangle - y})$ are both 1-q.s.c. losses.

For the remainder of the problem, assume that the data $\mathcal{X} \subset \mathbb{R}^d$ satisfy $\|x\|_2 \leq \sqrt{d}$ for all $x \in \mathcal{X}$. Let $L(\theta) = \mathbb{E}[\ell(\theta, (X, Y))]$ and $\theta^* = \operatorname{argmin}_{\theta} L(\theta)$. Assume that $\nabla^2 L(\theta^*) \succeq \lambda I$, where $\lambda > 0$ is fixed, and let $L_n(\theta) = P_n \ell(\theta, (X, Y))$ as usual.

(b) Show that if $\|\theta - \theta^*\|_2 \leq 1/\sqrt{d}$, then $\nabla^2 L_n(\theta) \succeq e^{-C} \nabla^2 L_n(\theta^*)$.

(c) Argue that if $t = \|\theta - \theta^*\|_2 \leq \frac{1}{\sqrt{d}}$ and $\|\nabla L_n(\theta^*)\|_2 \leq \gamma$, then

$$L_n(\theta) \geq L_n(\theta^*) - \gamma t + \frac{e^{-C} \lambda_{\min}(\nabla^2 L_n(\theta^*))}{2} t^2.$$

(d) Let ℓ be a 1-q.s.c. loss and assume that $|h'(t, y)| \leq 1$ and $|h''(t, y)| \leq 1$ for all $t \in \mathbb{R}$. Give a result similar to that of Theorem 5.3.8, but show that your conclusions hold with probability at least $1 - \delta$ as soon as

$$n \gtrsim \frac{d^2}{\lambda^2} \log \frac{d}{\delta}.$$

Exercise 5.7 (Truncation to obtain a moment bound): Let $B < \infty$. Show that under the conditions of Corollary 5.3.10,

$$\mathbb{E} \left[\|\hat{\theta}_n - \theta^*\|_2^2 \wedge B \right] \leq \frac{\text{tr}(\nabla^2 L(\theta^*)^{-1} \text{Cov}(\nabla \ell(\theta^*, Z)) \nabla^2 L(\theta^*)^{-1})}{n} + \frac{C \log n}{n^{3/2}},$$

where C is a problem-dependent constant.

Exercise 5.8: Let $\hat{\theta}_n = \text{argmin}_{\theta} L_n(\theta)$ for an M-estimation problem satisfying the conditions of Corollary 5.3.10. Show that for any $k \geq 1$, there are events \mathcal{E}_n with $\mathbb{P}(\mathcal{E}_n) \geq 1 - n^{-k}$ and for which

$$\left| \mathbb{E} \left[\|\hat{\theta}_n - \theta^*\|_2^2 \mathbf{1}_{\{\mathcal{E}_n\}} \right] - \frac{1}{n} \text{tr}(\nabla^2 L(\theta^*)^{-1} \text{Cov}(\nabla \ell(\theta^*, Z)) \nabla^2 L(\theta^*)^{-1}) \right| \leq \frac{Ck \log n}{n^{3/2}},$$

where C is a problem-dependent constant.

Exercise 5.9: In this problem, you provide sufficient conditions for exponential family models to have minimizers.

- (a) Let P_{θ} be a minimal exponential family (Definition 3.2) with density $p_{\theta}(x) = \exp(\theta^{\top} x - A(\theta))$ with respect to a base measure μ . Show that for any $\theta^* \in \text{dom } A$, the well-specified population loss $L(\theta) = -\mathbb{E}_{\theta^*}[\log p_{\theta}(X)]$ has unique minimizer θ^* .

For the remainder of the problem, we no longer assume the model is well-specified. For a measure μ on a set \mathcal{X} , recall that the essential supremum of a function $f : \mathcal{X} \rightarrow \mathbb{R}$ is

$$\text{ess sup}_{\mu}(f) := \inf \{t \in \mathbb{R} \mid \mu(\{f(x) \geq t\}) = 0\}.$$

We say that a measure μ on \mathbb{R}^d *essentially covers* a vector $v \in \mathbb{R}^d$ if

$$v^{\top} u < \text{ess sup}_{\mu} \{x^{\top} u\} \quad \text{for all } u \neq 0. \quad (5.4.1)$$

- (b) Let $\{P_{\theta}\}$ be the exponential family with density $p_{\theta}(x) = \exp(\theta^{\top} x - A(\theta))$ with respect to the measure μ . Let X be a random variable for which μ essentially covers (5.4.1) the mean $\mathbb{E}[X]$. Show that $L(\theta) = -\mathbb{E}[\log p_{\theta}(X)]$ has a minimizer.

Hint. A continuous convex function h has a minimizer if it is coercive, meaning that $h(\theta) \rightarrow \infty$ whenever $\|\theta\|_2 \rightarrow \infty$. Corollary C.3.7, part (i) may be useful.

Exercise 5.10: In this problem, you provide sufficient conditions for generalized linear models to have minimizers. Let $L(\theta) = \mathbb{E}[\ell(\theta, (X, Y))]$ be the population loss (which may be misspecified).

- (a) Consider Poisson regression with loss $\ell(\theta, (x, y)) = -yx^{\top} \theta + e^{\theta^{\top} x}$. Show that L has a unique minimizer if $\mathbb{E}[X] = 0$ and $\text{Cov}(X) = \mathbb{E}[XX^{\top}] \succ 0$.
- (b) Consider logistic regression with loss $\ell(\theta, (x, y)) = -yx^{\top} \theta + \log(1 + e^{\theta^{\top} x})$ for $y \in \{0, 1\}$. Show that L has a unique minimizer if $\mathbb{E}[XX^{\top}] \succ 0$ and $0 < \mathbb{P}(Y = 1 \mid X) < 1$ with probability 1 over X .

- (c) Consider a generalized linear model with densities $p_\theta(y | x) = \exp(\phi(x, y)^\top \theta - A(\theta | x))$ w.r.t. a base measure $\mu(\cdot | x)$ on $y \in \mathcal{Y}$, and assume for simplicity that $\mu(\mathcal{Y} | x) = 1$ for all x . Assume that for each vector $v \in \mathbb{S}^{d-1}$ and $x \in \mathcal{X}$,

$$\operatorname{ess\,sup}_{\mu(\cdot | x)} \{v^\top \phi(x, y)\} \geq \mathbb{E}_P[v^\top \phi(x, Y) | X = x],$$

and the set of x for which a strict inequality holds has positive P -probability. (This is equivalent to the set of x for which $\mu(\cdot | x)$ essentially covers (5.4.1) the conditional mean $\mathbb{E}_P[\phi(x, Y) | X = x]$ having positive probability.) Show that a minimizer of L exists. You may assume $\mathbb{E}[|A(\theta | X)|] < \infty$ for $\|\theta\|_2 \leq 1$ if it is convenient.

Hint. The techniques to solve Exercise 5.9 may be useful. In addition, see Exercise 4.2.

Exercise 5.11: Consider the robust regression setting of Example 5.3.3, and let $h \geq 0$ be a symmetric convex function, twice continuously differentiable in a neighborhood of 0. Assume that for any (measurable) subset $\mathcal{X}_0 \subset \mathcal{X}$, $\mathbb{E}[Y | X \in \mathcal{X}_0]$ exists and is finite, and assume $\mathbb{E}[XX^\top] \succ 0$. Show that a minimizer of $L(\theta) := \mathbb{E}[h(\langle \theta, X \rangle - Y)]$ exists. *Hint.* Show that L is coercive. Corollary C.3.7, part (i) may be useful.

Exercise 5.12 (The delta method for approximate sums): Let $T : \mathbb{R}^d \rightarrow \mathbb{R}^p$ be a differentiable function with derivative matrix $\dot{T}(\theta) \in \mathbb{R}^{p \times d}$, so that $T(\theta + \Delta) = T(\theta) + \dot{T}(\theta)\Delta + o(\|\Delta\|)$ as $\Delta \rightarrow 0$. Let $\hat{\theta}_n \in \mathbb{R}^d$ be a sequence of random vectors with

$$\hat{\theta}_n - \theta = P_n Z + R_n,$$

where Z_i are i.i.d. and R_n is a remainder term.

- (a) Assume that $\mathbb{E}[\|Z_i\|_2^2] < \infty$ and that for each $\epsilon > 0$, $\mathbb{P}(\|R_n\|_2 \geq \epsilon/\sqrt{n}) \rightarrow 0$ as $n \rightarrow \infty$. Show that

$$\hat{\theta}_n - \theta = \dot{T}(\theta)P_n Z + R'_n,$$

where the remainder R'_n also satisfies $\mathbb{P}(\|R'_n\|_2 \geq \epsilon/\sqrt{n}) \rightarrow 0$ for all $\epsilon > 0$.

- (b) Assume that T is locally smooth enough that for some $K < \infty$ and $\delta > 0$,

$$\left\|T(\theta + \Delta) - T(\theta) - \dot{T}(\theta)\Delta\right\|_2 \leq K \|\Delta\|_2^2$$

when $\|\Delta\|_2 \leq \delta$. Assume additionally that there exist $C_0, C_1 < \infty$ such that for $t \geq 0$, we have $\|R_n\|_2 \leq \frac{C_0 t}{n}$ with probability at least $1 - e^{-t}$ and that $\mathbb{P}(\|P_n Z\|_2 \geq t) \leq C_1 \exp(-nt^2/\sigma^2)$. Give a quantitative version of part (a).

JCD Comment: Add in some connections to the exponential family material. Some ideas:

1. A hypothesis test likelihood ratio for them (see page 40 of handwritten notes)
2. A full learning guarantee with convergence of Hessian and everything, e.g., for logistic regression?
3. In the Ledoux-Talagrand stuff, maybe worth going through example of logistic regression. Also, having working logistic example throughout? Helps clear up the structure and connect with exponential families.
4. Maybe an exercise for Lipschitz functions with random Lipschitz constants?

Chapter 6

Generalization and stability

Concentration inequalities provide powerful techniques for demonstrating when random objects that are functions of collections of independent random variables—whether sample means, functions with bounded variation, or collections of random vectors—behave similarly to their expectations. This chapter continues exploration of these ideas by incorporating the central thesis of this book: that information theory’s connections to statistics center around measuring when (and how) two probability distributions get close to one another. On its face, we remain focused on the main objects of the preceding chapter, where we have a population probability distribution P on a space \mathcal{X} and some collection of functions $f : \mathcal{X} \rightarrow \mathbb{R}$. We then wish to understand when we expect the empirical distribution

$$P_n := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i},$$

defined by the sample $X_i \stackrel{\text{iid}}{\sim} P$, to be close to the population P as measured by f . Following the notation we introduce in Section 5.1, for $Pf := \mathbb{E}_P[f(X)]$, we again ask to have

$$P_n f - Pf = \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E}_P[f(X)])$$

to be small simultaneously for all f .

In this chapter, however, we develop a family of tools based around *PAC* (*probably approximately correct*) *Bayesian* bounds, where we slightly perturb the functions f of interest to average them in some way; when these perturbations keep $P_n f$ stable, we expect that $P_n f \approx Pf$, that is, the sample generalizes to the population. These perturbations allow us to bring the tools of the divergence measures we have developed to bear on the problems of convergence and generalization. They also allow us to go beyond the “basic” concentration inequalities to situations with interaction, where a data analyst may evaluate some functions of P_n , then adaptively choose additional queries or analyses to do on the sample X_1^n . This breaks standard statistical analyses—which assume an *a priori* specified set of hypotheses or questions to be answered—but is possible to address once we can limit the information the analyses release in precise ways that information-theoretic tools allow. Even more, in the next chapter we show how they form the basis for *transportation inequalities*, powerful tools for concentration of measure. Modern work has also shown how to leverage these techniques, coupled with computation, to provide non-vacuous bounds on learning for complicated scenarios and models to which all classical bounds fail to apply, such as deep learning.

6.1 The variational representation of Kullback-Leibler divergence

The starting point of all of our generalization bounds is a surprisingly simply variational result, which relates expectations, moment generating functions, and the KL-divergence in one single equality. It turns out that this inequality, by relating means with moment generating functions and divergences, allows us to prove generalization bounds based on information-theoretic tools and stability.

Theorem 6.1.1 (Donsker-Varadhan variational representation). *Let P and Q be distributions on a common space \mathcal{X} . Then*

$$D_{\text{kl}}(P\|Q) = \sup_g \left\{ \mathbb{E}_P[g(X)] - \log \mathbb{E}_Q[e^{g(X)}] \right\},$$

where the supremum is taken over measurable functions $g : \mathcal{X} \rightarrow \mathbb{R}$ with $\mathbb{E}_Q[e^{g(X)}] < \infty$. We can also replace this by bounded simple functions g .

We give one proof of this result and one sketch of a proof, which holds when the underlying space is discrete, that may be more intuitive: the first constructs a particular “tilting” of Q via the function e^g , and verifies the equality. The second relies on the discretization of the KL-divergence and may be more intuitive to readers familiar with convex optimization: essentially, we expect this result because the function $\log(\sum_{j=1}^k e^{x_j})$ is the convex conjugate of the negative entropy. (See also Exercise 6.1.)

Proof We may assume that P is absolutely continuous with respect to Q , meaning that $Q(A) = 0$ implies that $P(A) = 0$, as otherwise both sides are infinite by inspection. Thus, it is no loss of generality to let P and Q have densities p and q .

Attainment in the equality is easy: we simply take $g(x) = \log \frac{p(x)}{q(x)}$, so that $\mathbb{E}_Q[e^{g(X)}] = 1$. To show that the right hand side is never larger than $D_{\text{kl}}(P\|Q)$ requires a bit more work. To that end, let g be any function such that $\mathbb{E}_Q[e^{g(X)}] < \infty$, and define the random variable $Z_g(x) = e^{g(x)}/\mathbb{E}_Q[e^{g(X)}]$, so that $\mathbb{E}_Q[Z_g] = 1$. Then using the absolute continuity of P w.r.t. Q , we have

$$\begin{aligned} \mathbb{E}_P[\log Z_g] &= \mathbb{E}_P \left[\log \frac{p(X)}{q(X)} + \log \left(Z_g(X) \frac{q(X)}{p(X)} \right) \right] = D_{\text{kl}}(P\|Q) + \mathbb{E}_P \left[\log \left(Z_g \frac{dQ}{dP} \right) \right] \\ &\leq D_{\text{kl}}(P\|Q) + \log \mathbb{E}_P \left[\frac{dQ}{dP} Z_g \right] \\ &= D_{\text{kl}}(P\|Q) + \log \mathbb{E}_Q[Z_g]. \end{aligned}$$

As $\mathbb{E}_Q[Z_g] = 1$, using that $\mathbb{E}_P[\log Z_g] = \mathbb{E}_P[g(X)] - \log \mathbb{E}_Q[e^{g(X)}]$ gives the result.

For the claim that bounded simple functions are sufficient, all we need to do is demonstrate (asymptotic) achievability. For this, we use the definition (2.2.1) of the KL-divergence as a supremum over partitions. Take \mathcal{A}_n be an sequence of partitions so that $D_{\text{kl}}(P\|Q | \mathcal{A}_n) \rightarrow D_{\text{kl}}(P\|Q)$. Then let $g_n(x) = \sum_{A \in \mathcal{A}_n} \mathbf{1}\{x \in A\} \log \frac{P(A)}{Q(A)}$, which gives $D_{\text{kl}}(P\|Q | \mathcal{A}_n) = \mathbb{E}_P[g_n(X)] - \log \mathbb{E}_Q[e^{g_n(X)}]$. \square

Here is the second proof of Theorem 6.1.1, which applies when \mathcal{X} is discrete and finite. That we can approximate KL-divergence by suprema over finite partitions (as in definition (2.2.1)) suggests that this approach works in general—which it can—but this requires some not completely trivial

approximations of $\mathbb{E}_P[g]$ and $\mathbb{E}_Q[e^g]$ by discretized versions of their expectations, which makes things rather tedious.

Proof of Theorem 6.1.1, the finite case As we have assumed that P and Q have finite supports, which we identify with $\{1, \dots, k\}$ and p.m.f.s $p, q \in \Delta_k = \{p \in \mathbb{R}_+^k \mid \langle \mathbf{1}, p \rangle = 1\}$. Define $f_q(v) = \log(\sum_{j=1}^k q_j e^{v_j})$, which is convex in v (recall Proposition 3.2.1). Then the supremum in the variational representation takes the form

$$h(p) := \sup_{v \in \mathbb{R}^k} \{\langle p, v \rangle - f_q(v)\}.$$

If we can take derivatives and solve for zero, we are guaranteed to achieve the supremum. To that end, note that

$$\nabla_v \{\langle p, v \rangle - f_q(v)\} = p - \left[\frac{q_i e^{v_i}}{\sum_{j=1}^k q_j e^{v_j}} \right]_{i=1}^k,$$

so that setting $v_j = \log \frac{p_j}{q_j}$ achieves $p - \nabla_v f_q(v) = p - p = 0$ and hence the supremum. Noting that $\log(\sum_{j=1}^k q_j \exp(\log \frac{p_j}{q_j})) = \log(\sum_{j=1}^k p_j) = 0$ gives $h(p) = D_{\text{kl}}(p \| q)$. \square

The Donsker-Varadhan variational representation already gives a hint that we can use some information-theoretic techniques to control the difference between an empirical sample and its expectation, at least in an average sense. In particular, we see that for any function g , we have

$$\mathbb{E}_P[g(X)] \leq D_{\text{kl}}(P \| Q) + \log \mathbb{E}_Q[e^{g(X)}]$$

for any random variable X . Now, changing this on its head a bit, suppose that we consider a collection of functions \mathcal{F} and put two probability measures π and π_0 on \mathcal{F} , and consider $P_n f - P f$, where we consider f a random variable $f \sim \pi$ or $f \sim \pi_0$. Then a consequence of the Donsker-Varadhan theorem is that

$$\int (P_n f - P f) d\pi(f) \leq D_{\text{kl}}(\pi \| \pi_0) + \log \int \exp(P_n f - P f) d\pi_0(f)$$

for any π, π_0 . While this inequality is a bit naive—bounding a difference by an exponent seems wasteful—as we shall see, it has substantial applications when we can upper bound the KL-divergence $D_{\text{kl}}(\pi \| \pi_0)$.

6.2 PAC-Bayes bounds

Probably-approximately-correct (PAC) Bayesian bounds proceed from a perspective similar to that of the covering numbers and covering entropies we develop in Section 5.1, where if for a collection of functions \mathcal{F} there is a finite subset (a cover) $\{f_v\}$ such that each $f \in \mathcal{F}$ is “near” one of the f_v , then we need only control deviations of $P_n f$ from $P f$ for the elements of $\{f_v\}$. In PAC-Bayes bounds, we instead average functions f with other functions, and this averaging allows a similar family of guarantees and applications.

Let us proceed with the main results. Let \mathcal{F} be a collection of functions $f : \mathcal{X} \rightarrow \mathbb{R}$, and assume that each function f is σ^2 -sub-Gaussian, which we recall (Definition 4.1) means that $\mathbb{E}[e^{\lambda(f(X) - P f)}] \leq \exp(\lambda^2 \sigma^2 / 2)$ for all $\lambda \in \mathbb{R}$, where $P f = \mathbb{E}_P[f(X)] = \int f(x) dP(x)$ denotes the

expectation of f under P . The main theorem of this section shows that averages of the squared error $(P_n f - P f)^2$ of the empirical distribution P_n to P converge quickly to zero for *all* averaging distributions π on functions $f \in \mathcal{F}$ so long as each f is σ^2 -sub-Gaussian, with the caveat that we pay a cost for different choices of π . The key is that we choose some prior distribution π_0 on \mathcal{F} first.

Theorem 6.2.1. *Let Π be the collection of all probability distributions on the set \mathcal{F} and let π_0 be a fixed prior probability distribution on $f \in \mathcal{F}$. With probability at least $1 - \delta$,*

$$\int (P_n f - P f)^2 d\pi(f) \leq \frac{8\sigma^2}{3} \frac{D_{\text{kl}}(\pi \| \pi_0) + \log \frac{2}{\delta}}{n} \quad \text{simultaneously for all } \pi \in \Pi.$$

Proof The key is to combine Example 4.1.12 with the variational representation that Theorem 6.1.1 provides for KL-divergences. We state Example 4.1.12 as a lemma here.

Lemma 6.2.2. *Let Z be a σ^2 -sub-Gaussian random variable. Then for $\lambda \geq 0$,*

$$\mathbb{E}[e^{\lambda Z^2}] \leq \frac{1}{\sqrt{[1 - 2\sigma^2\lambda]_+}}.$$

Without loss of generality, we assume that $P f = 0$ for all $f \in \mathcal{F}$, and recall that $P_n f = \frac{1}{n} \sum_{i=1}^n f(X_i)$ is the empirical mean of f . Then we know that $P_n f$ is σ^2/n -sub-Gaussian, and Lemma 6.2.2 implies that $\mathbb{E}[\exp(\lambda(P_n f)^2)] \leq [1 - 2\lambda\sigma^2/n]_+^{-1/2}$ for any f , and thus for any prior π_0 on f we have

$$\mathbb{E} \left[\int \exp(\lambda(P_n f)^2) d\pi_0(f) \right] \leq [1 - 2\lambda\sigma^2/n]_+^{-1/2}.$$

Consequently, taking $\lambda = \lambda_n := \frac{3n}{8\sigma^2}$, we obtain

$$\mathbb{E} \left[\int \exp(\lambda_n(P_n f)^2) d\pi_0(f) \right] = \mathbb{E} \left[\int \exp \left(\frac{3n}{8\sigma^2} (P_n f)^2 \right) d\pi_0(f) \right] \leq 2.$$

Markov's inequality thus implies that

$$\mathbb{P} \left(\int \exp(\lambda_n(P_n f)^2) d\pi_0(f) \geq \frac{2}{\delta} \right) \leq \delta, \quad (6.2.1)$$

where the probability is over $X_i \stackrel{\text{iid}}{\sim} P$.

Now, we use the Donsker-Varadhan equality (Theorem 6.1.1). Letting $\lambda > 0$, we define the function $g(f) = \lambda(P_n f)^2$, so that for any two distributions π and π_0 on \mathcal{F} , we have

$$\frac{1}{\lambda} \int g(f) d\pi(f) = \int (P_n f)^2 d\pi(f) \leq \frac{D_{\text{kl}}(\pi \| \pi_0) + \log \int \exp(\lambda(P_n f)^2) d\pi_0(f)}{\lambda}.$$

This holds without any probabilistic qualifications, so using the application (6.2.1) of Markov's inequality with $\lambda = \lambda_n$, we thus see that with probability at least $1 - \delta$ over X_1, \dots, X_n , simultaneously for all distributions π ,

$$\int (P_n f)^2 d\pi(f) \leq \frac{8\sigma^2}{3} \frac{D_{\text{kl}}(\pi \| \pi_0) + \log \frac{2}{\delta}}{n}.$$

This is the desired result (as we have assumed that $Pf = 0$ w.l.o.g.). \square

By Jensen's inequality (or Cauchy-Schwarz), it is immediate from Theorem 6.2.1 that we also have

$$\int |P_n f - P f| d\pi(f) \leq \sqrt{\frac{8\sigma^2}{3} \frac{D_{\text{kl}}(\pi \| \pi_0) + \log \frac{2}{\delta}}{n}} \quad \text{simultaneously for all } \pi \in \Pi \quad (6.2.2)$$

with probability at least $1 - \delta$, so that $\mathbb{E}_\pi[|P_n f - P f|]$ is with high probability of order $1/\sqrt{n}$. The inequality (6.2.2) is the original form of the PAC-Bayes bound due to McAllester, with slightly sharper constants and improved logarithmic dependence. The key is that *stability*, in the form of a prior π_0 and posterior π closeness, allow us to achieve reasonably tight control over the deviations of random variables and functions with high probability.

Let us give an example, which is similar to many of our approaches in Section 5.2, to illustrate some of the approaches this allows. The basic idea is that by appropriate choice of prior π_0 and “posterior” π , whenever we have appropriately smooth classes of functions we achieve certain generalization guarantees.

Example 6.2.3 (A uniform law for Lipschitz functions): Consider a case as in Section 5.2, where we let $L(\theta) = P\ell(\theta, Z)$ for some function $\ell : \Theta \times \mathcal{Z} \rightarrow \mathbb{R}$. Let $\mathbb{B}_2^d = \{v \in \mathbb{R}^d \mid \|v\|_2 \leq 1\}$ be the ℓ_2 -ball in \mathbb{R}^d , and let us assume that $\Theta \subset r\mathbb{B}_2^d$ and additionally that $\theta \mapsto \ell(\theta, z)$ is M -Lipschitz for all $z \in \mathcal{Z}$. For simplicity, we assume that $\ell(\theta, z) \in [0, 2Mr]$ for all $\theta \in \Theta$ (we may simply relativize our bounds by replacing ℓ by $\ell(\cdot, z) - \inf_{\theta \in \Theta} \ell(\theta, z) \in [0, 2Mr]$).

If $\hat{L}_n(\theta) = P_n \ell(\theta, Z)$, then Theorem 6.2.1 implies that

$$\int |\hat{L}_n(\theta) - L(\theta)| d\pi(\theta) \leq \sqrt{\frac{8M^2 r^2}{3n} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{2}{\delta} \right]}$$

for all π with probability at least $1 - \delta$. Now, let $\theta_0 \in \Theta$ be arbitrary, and for $\epsilon > 0$ (to be chosen later) take π_0 to be uniform on $(r + \epsilon)\mathbb{B}_2^d$ and π to be uniform on $\theta_0 + \epsilon\mathbb{B}_2^d$. Then we immediately see that $D_{\text{kl}}(\pi \| \pi_0) = d \log(1 + \frac{r}{\epsilon})$. Moreover, we have $\int \hat{L}_n(\theta) d\pi(\theta) \in \hat{L}_n(\theta_0) \pm M\epsilon$ and similarly for $L(\theta)$, by the M -Lipschitz continuity of ℓ . For any fixed $\epsilon > 0$, we thus have

$$|\hat{L}_n(\theta_0) - L(\theta_0)| \leq 2M\epsilon + \sqrt{\frac{2M^2 r^2}{3n} \left[d \log \left(1 + \frac{r}{\epsilon} \right) + \log \frac{2}{\delta} \right]}$$

simultaneously for all $\theta_0 \in \Theta$, with probability at least $1 - \delta$. By choosing $\epsilon = \frac{rd}{n}$ we obtain that with probability at least $1 - \delta$,

$$\sup_{\theta \in \Theta} |\hat{L}_n(\theta) - L(\theta)| \leq \frac{2Mr d}{n} + \sqrt{\frac{8M^2 r^2}{3n} \left[d \log \left(1 + \frac{n}{d} \right) + \log \frac{2}{\delta} \right]}.$$

Thus, roughly, with high probability we have $|\hat{L}_n(\theta) - L(\theta)| \leq O(1)Mr\sqrt{\frac{d}{n} \log \frac{n}{d}}$ for all θ . \diamond

On the one hand, the result in Example 6.2.3 is satisfying: it applies to any Lipschitz function and provides a uniform bound. On the other hand, when we compare to the results achievable for

specially structured linear function classes, then applying Rademacher complexity bounds—such as Proposition 5.2.9 and Example 5.2.10—we have somewhat weaker results, in that they depend on the dimension explicitly, while the Rademacher bounds do not exhibit this explicit dependence. This means they can potentially apply in infinite dimensional spaces that Example 6.2.3 cannot. We will give an example presently showing how to address some of these issues.

6.2.1 Relative bounds

In many cases, it is useful to have bounds that provide somewhat finer control than the bounds we have presented. Recall from our discussion of sub-Gaussian and sub-exponential random variables, especially the Bennett and Bernstein-type inequalities (Proposition 4.1.21), that if a random variable X satisfies $|X| \leq b$ but $\text{Var}(X) \leq \sigma^2 \ll b^2$, then X concentrates more quickly about its mean than the convergence provided by naive application of sub-Gaussian concentration with sub-Gaussian parameter $b^2/8$. To that end, we investigate an alternative to Theorem 6.2.1 that allows somewhat sharper control.

The approach is similar to our derivation in Theorem 6.2.1, where we show that the moment generating function of a quantity like $P_n f - P f$ is small (Eq. (6.2.1)) and then relate this—via the Donsker-Varadhan change of measure in Theorem 6.1.1—to the quantities we wish to control. In the next proposition, we provide relative bounds on the deviations of functions from their means. To make this precise, let \mathcal{F} be a collection of functions $f : \mathcal{X} \rightarrow \mathbb{R}$, and let $\sigma^2(f) := \text{Var}(f(X))$ be the variance of functions in \mathcal{F} . We assume the class satisfies the Bernstein condition (4.1.7) with parameter b , that is,

$$\left| \mathbb{E} \left[(f(X) - P f)^k \right] \right| \leq \frac{k!}{2} \sigma^2(f) b^{k-2} \quad \text{for } k = 3, 4, \dots \quad (6.2.3)$$

This says that the second moment of functions $f \in \mathcal{F}$ bounds—with the additional boundedness-type constant b —the higher moments of functions in \mathcal{F} . We then have the following result.

Proposition 6.2.4. *Let \mathcal{F} be a collection of functions $f : \mathcal{X} \rightarrow \mathbb{R}$ satisfying the Bernstein condition (6.2.3). Then for any $|\lambda| \leq \frac{1}{2b}$, with probability at least $1 - \delta$,*

$$\lambda \int P f d\pi(f) - \lambda^2 \int \sigma^2(f) d\pi(f) \leq \lambda \int P_n f d\pi(f) + \frac{1}{n} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{1}{\delta} \right]$$

simultaneously for all $\pi \in \Pi$.

Proof We begin with an inequality on the moment generating function of random variables satisfying the Bernstein condition (4.1.7), that is, that $|\mathbb{E}[(X - \mu)^k]| \leq \frac{k!}{2} \sigma^2 b^{k-2}$ for $k \geq 2$. In this case, Lemma 4.1.20 implies that

$$\mathbb{E}[e^{\lambda(X - \mu)}] \leq \exp(\lambda^2 \sigma^2)$$

for $|\lambda| \leq 1/(2b)$. As a consequence, for any f in our collection \mathcal{F} , we see that if we define

$$\Delta_n(f, \lambda) := \lambda [P_n f - P f - \lambda \sigma^2(f)],$$

we have that

$$\mathbb{E}[\exp(n \Delta_n(f, \lambda))] = \mathbb{E}[\exp(\lambda(f(X) - P f) - \lambda^2 \sigma^2(f))]^n \leq 1$$

for all n , $f \in \mathcal{F}$, and $|\lambda| \leq \frac{1}{2b}$. Then, for any fixed measure π_0 on \mathcal{F} , Markov's inequality implies that

$$\mathbb{P} \left(\int \exp(n \Delta_n(f, \lambda)) d\pi_0(f) \geq \frac{1}{\delta} \right) \leq \delta. \quad (6.2.4)$$

Now, as in the proof of Theorem 6.2.1, we use the Donsker-Varadhan Theorem 6.1.1 (change of measure), which implies that

$$n \int \Delta_n(f, \lambda) d\pi(f) \leq D_{\text{kl}}(\pi \| \pi_0) + \log \int \exp(n \Delta_n(f, \lambda)) d\pi_0(f)$$

for all distributions π . Using inequality (6.2.4), we obtain that with probability at least $1 - \delta$,

$$\int \Delta_n(f, \lambda) d\pi(f) \leq \frac{1}{n} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{1}{\delta} \right]$$

for all π . As this holds for any fixed $|\lambda| \leq 1/(2b)$, this gives the desired result by rearranging. \square

We would like to optimize over the bound in Proposition 6.2.4 by choosing the “best” λ . If we *could* choose the optimal λ , by rearranging Proposition 6.2.4 we would obtain the bound

$$\begin{aligned} \mathbb{E}_\pi[Pf] &\leq \mathbb{E}_\pi[P_n f] + \inf_{\lambda > 0} \left\{ \lambda \mathbb{E}_\pi[\sigma^2(f)] + \frac{1}{n\lambda} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{1}{\delta} \right] \right\} \\ &= \mathbb{E}_\pi[P_n f] + 2 \sqrt{\frac{\mathbb{E}_\pi[\sigma^2(f)]}{n} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{1}{\delta} \right]} \end{aligned}$$

simultaneously for all π , with probability at least $1 - \delta$. The problem with this approach is two-fold: first, we cannot arbitrarily choose λ in Proposition 6.2.4, and second, the bound above depends on the unknown population variance $\sigma^2(f)$. It is thus of interest to understand situations in which we can obtain similar guarantees, but where we can replace unknown population quantities on the right side of the bound with known quantities.

To that end, let us consider the following condition, a type of relative error condition related to the Bernstein condition (4.1.7): for each $f \in \mathcal{F}$,

$$\sigma^2(f) \leq bPf. \tag{6.2.5}$$

This condition is most natural when each of the functions f take nonnegative values—for example, when $f(X) = \ell(\theta, X)$ for some loss function ℓ and parameter θ of a model. If the functions f are nonnegative and upper bounded by b , then we certainly have $\sigma^2(f) \leq \mathbb{E}[f(X)^2] \leq b\mathbb{E}[f(X)] = bPf$, so that Condition (6.2.5) holds. Revisiting Proposition 6.2.4, we rearrange to obtain the following theorem.

Theorem 6.2.5. *Let \mathcal{F} be a collection of functions satisfying the Bernstein condition (6.2.3) as in Proposition 6.2.4, and in addition, assume the variance-bounding condition (6.2.5). Then for any $0 \leq \lambda \leq \frac{1}{2b}$, with probability at least $1 - \delta$,*

$$\mathbb{E}_\pi[Pf] \leq \mathbb{E}_\pi[P_n f] + \frac{\lambda b}{1 - \lambda b} \mathbb{E}_\pi[P_n f] + \frac{1}{\lambda(1 - \lambda b)} \frac{1}{n} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{1}{\delta} \right]$$

for all π .

Proof We use condition (6.2.5) to see that

$$\lambda \mathbb{E}_\pi[Pf] - \lambda^2 b \mathbb{E}_\pi[Pf] \leq \lambda \mathbb{E}_\pi[Pf] - \lambda^2 \mathbb{E}_\pi[\sigma^2(f)],$$

apply Proposition 6.2.4, and divide both sides of the resulting inequality by $\lambda(1 - \lambda b)$. \square

To make this uniform in λ , thus achieving a tighter bound (so that we need not pre-select λ), we choose multiple values of λ and apply a union bound. To that end, let $1 + \eta = \frac{1}{1 - \lambda b}$, or $\eta = \frac{\lambda b}{1 - \lambda b}$ and $\frac{1}{\lambda b(1 - \lambda b)} = \frac{(1 + \eta)^2}{\eta}$, so that the inequality in Theorem 6.2.1 is equivalent to

$$\mathbb{E}_\pi[Pf] \leq \mathbb{E}_\pi[P_nf] + \eta \mathbb{E}_\pi[P_nf] + \frac{(1 + \eta)^2}{\eta} \frac{b}{n} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{1}{\delta} \right].$$

Using that our choice of $\eta \in [0, 1]$, this implies

$$\mathbb{E}_\pi[Pf] \leq \mathbb{E}_\pi[P_nf] + \eta \mathbb{E}_\pi[P_nf] + \frac{1}{\eta} \frac{b}{n} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{1}{\delta} \right] + \frac{3b}{n} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{1}{\delta} \right].$$

Now, take $\eta_1 = 1/n, \dots, \eta_n = 1$. Then by optimizing over $\eta \in \{\eta_1, \dots, \eta_n\}$ (which is equivalent, to within a $1/n$ factor, to optimizing over $0 < \eta \leq 1$) and applying a union bound, we obtain

Corollary 6.2.6. *Let the conditions of Theorem 6.2.5 hold. Then with probability at least $1 - \delta$,*

$$\mathbb{E}_\pi[Pf] \leq \mathbb{E}_\pi[P_nf] + 2\sqrt{\frac{b\mathbb{E}_\pi[P_nf]}{n}} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{n}{\delta} \right] + \frac{1}{n} \left(\mathbb{E}_\pi[P_nf] + 5b \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{n}{\delta} \right] \right),$$

simultaneously for all π on \mathcal{F} .

Proof By a union bound, we have

$$\mathbb{E}_\pi[Pf] \leq \mathbb{E}_\pi[P_nf] + \eta \mathbb{E}_\pi[P_nf] + \frac{1}{\eta} \frac{b}{n} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{n}{\delta} \right] + \frac{3b}{n} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{n}{\delta} \right]$$

for each $\eta \in \{1/n, \dots, 1\}$. We consider two cases. In the first, assume that $\mathbb{E}_\pi[P_nf] \leq \frac{b}{n} (D_{\text{kl}}(\pi \| \pi_0) + \log \frac{n}{\delta})$. Then taking $\eta = 1$ above evidently gives the result. In the second, we have $\mathbb{E}_\pi[P_nf] > \frac{b}{n} (D_{\text{kl}}(\pi \| \pi_0) + \log \frac{n}{\delta})$, and we can set

$$\eta_\star = \sqrt{\frac{\frac{b}{n} (D_{\text{kl}}(\pi \| \pi_0) + \log \frac{n}{\delta})}{\mathbb{E}_\pi[P_nf]}} \in (0, 1).$$

Choosing η to be the smallest value η_k in $\{\eta_1, \dots, \eta_n\}$ with $\eta_k \geq \eta_\star$, so that $\eta_\star \leq \eta \leq \eta_\star + \frac{1}{n}$ then implies the claim in the corollary. \square

6.2.2 A large-margin guarantee

Let us revisit the loss minimization approaches central to Section 5.2 and Example 6.2.3 in the context of Corollary 6.2.6. We will investigate an approach to achieve convergence guarantees that are (nearly) independent of dimension, focusing on 0-1 losses in a binary classification problem. Consider a binary classification problem with data $(x, y) \in \mathbb{R}^d \times \{\pm 1\}$, where we make predictions $\langle \theta, x \rangle$ (or its sign), and for a *margin penalty* $\gamma \geq 0$ we define the loss

$$\ell_\gamma(\theta; (x, y)) = \mathbf{1} \{ \langle \theta, x \rangle y \leq \gamma \}.$$

We call the quantity $\langle \theta, x \rangle y$ the *margin* of θ on the pair (x, y) , noting that when the margin is large, $\langle \theta, x \rangle$ has the same sign as y and is “confident” (i.e. far from zero). For shorthand, let us define the expected and empirical losses at margin γ by

$$L_\gamma(\theta) := P\ell_\gamma(\theta; (X, Y)) \quad \text{and} \quad \widehat{L}_\gamma(\theta) := P_n\ell_\gamma(\theta; (X, Y)).$$

Consider the following scenario: the data x lie in a ball of radius b , so that $\|x\|_2 \leq b$; note that the losses ℓ_γ and ℓ_0 satisfy the Bernstein (6.2.3) and self-bounding (6.2.5) conditions with constant 1 as they take values in $\{0, 1\}$. We then have the following proposition.

Proposition 6.2.7. *Let the above conditions on the data (x, y) hold and let the margin $\gamma > 0$ and radius $r < \infty$. Then with probability at least $1 - \delta$,*

$$P(\langle \theta, X \rangle Y \leq 0) \leq \left(1 + \frac{1}{n}\right) P_n(\langle \theta, X \rangle Y \leq \gamma) + \sqrt{8 \frac{rb \log \frac{n}{\delta}}{\gamma \sqrt{n}}} \sqrt{P_n(\langle \theta, X \rangle Y \leq \gamma)} + C \frac{r^2 b^2 \log \frac{n}{\delta}}{\gamma^2 n}$$

simultaneously for all $\|\theta\|_2 \leq r$, where C is a numerical constant independent of the problem parameters.

Proposition 6.2.7 provides a “dimension-free” guarantee—it depends only on the ℓ_2 -norms $\|\theta\|_2$ and $\|x\|_2$ —so that it can apply equally in infinite dimensional spaces. The key to the inequality is that if we can find a large margin predictor—for example, one achieved by a support vector machine or, more broadly, by minimizing a convex loss of the form

$$\underset{\|\theta\|_2 \leq r}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \phi(\langle X_i, \theta \rangle Y_i)$$

for some decreasing convex $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$, e.g. $\phi(t) = [1 - t]_+$ or $\phi(t) = \log(1 + e^{-t})$ —then we get strong generalization performance guarantees relative to the empirical margin γ . As one particular instantiation of this approach, suppose we can obtain a perfect classifier with positive margin: a vector θ with $\|\theta\|_2 \leq r$ such that $\langle \theta, X_i \rangle Y_i \geq \gamma$ for each $i = 1, \dots, n$. Then Proposition 6.2.7 guarantees that

$$P(\langle \theta, X \rangle Y \leq 0) \leq C \frac{r^2 b^2 \log \frac{n}{\delta}}{\gamma^2 n}$$

with probability at least $1 - \delta$.

Proof Let π_0 be $\mathbf{N}(0, \tau^2 I)$ for some $\tau > 0$ to be chosen, and let π be $\mathbf{N}(\widehat{\theta}, \tau^2 I)$ for some $\widehat{\theta} \in \mathbb{R}^d$ satisfying $\|\widehat{\theta}\|_2 \leq r$. Then Corollary 6.2.6 implies that

$$\begin{aligned} & \mathbb{E}_\pi[L_\gamma(\theta)] \\ & \leq \mathbb{E}_\pi[\widehat{L}_\gamma(\theta)] + 2\sqrt{\frac{\mathbb{E}_\pi[\widehat{L}_\gamma(\theta)]}{n} \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{n}{\delta} \right]} + \frac{1}{n} \left(\mathbb{E}_\pi[\widehat{L}_\gamma(\theta)] + C \left[D_{\text{kl}}(\pi \| \pi_0) + \log \frac{n}{\delta} \right] \right) \\ & \leq \mathbb{E}_\pi[\widehat{L}_\gamma(\theta)] + 2\sqrt{\frac{\mathbb{E}_\pi[\widehat{L}_\gamma(\theta)]}{n} \left[\frac{r^2}{2\tau^2} + \log \frac{n}{\delta} \right]} + \frac{1}{n} \left(\mathbb{E}_\pi[\widehat{L}_\gamma(\theta)] + C \left[\frac{r^2}{2\tau^2} + \log \frac{n}{\delta} \right] \right) \end{aligned}$$

simultaneously for all $\widehat{\theta}$ satisfying $\|\widehat{\theta}\|_2 \leq r$ with probability at least $1 - \delta$, where we have used that $D_{\text{kl}}(\mathbf{N}(\widehat{\theta}, \tau^2 I) \| \mathbf{N}(0, \tau^2 I)) = \|\widehat{\theta}\|_2^2 / (2\tau^2)$.

Let us use the margin assumption. Note that if $Z \sim \mathbf{N}(0, \tau^2 I)$, then for any fixed θ_0, x, y we have

$$\ell_0(\theta_0; (x, y)) - \mathbb{P}(Z^\top x \geq \gamma) \leq \mathbb{E}[\ell_\gamma(\theta_0 + Z; (x, y))] \leq \ell_{2\gamma}(\theta_0; (x, y)) + \mathbb{P}(Z^\top x \geq \gamma)$$

where the middle expectation is over $Z \sim \mathbf{N}(0, \tau^2 I)$. Using the $\tau^2 \|x\|_2^2$ -sub-Gaussianity of $Z^\top x$, we can obtain immediately that if $\|x\|_2 \leq b$, we have

$$\ell_0(\theta_0; (x, y)) - \exp\left(-\frac{\gamma^2}{2\tau^2 b^2}\right) \leq \mathbb{E}[\ell_\gamma(\theta_0 + Z; (x, y))] \leq \ell_{2\gamma}(\theta_0; (x, y)) + \exp\left(-\frac{\gamma^2}{2\tau^2 b^2}\right).$$

Returning to our earlier bound, we evidently have that if $\|x\|_2 \leq b$ for all $x \in \mathcal{X}$, then with probability at least $1 - \delta$, simultaneously for all $\theta \in \mathbb{R}^d$ with $\|\theta\|_2 \leq r$,

$$\begin{aligned} L_0(\theta) &\leq \widehat{L}_{2\gamma}(\theta) + 2 \exp\left(-\frac{\gamma^2}{2\tau^2 b^2}\right) + 2 \sqrt{\frac{\widehat{L}_{2\gamma}(\theta) + \exp(-\frac{\gamma^2}{2\tau^2 b^2})}{n} \left[\frac{r^2}{2\tau^2} + \log \frac{n}{\delta}\right]} \\ &\quad + \frac{1}{n} \left(\widehat{L}_{2\gamma}(\theta) + \exp\left(-\frac{\gamma^2}{2\tau^2 b^2}\right) + C \left[\frac{r^2}{2\tau^2} + \log \frac{n}{\delta}\right] \right). \end{aligned}$$

Setting $\tau^2 = \frac{\gamma^2}{2b^2 \log n}$, we immediately see that for any choice of margin $\gamma > 0$, we have with probability at least $1 - \delta$ that

$$\begin{aligned} L_0(\theta) &\leq \widehat{L}_{2\gamma}(\theta) + \frac{2b}{n} + 2 \sqrt{\frac{1}{n} \left[\widehat{L}_{2\gamma}(\theta) + \frac{b}{n} \right] \left[\frac{r^2 b^2 \log n}{2\gamma^2} + \log \frac{n}{\delta} \right]} \\ &\quad + \frac{1}{n} \left(\widehat{L}_{2\gamma}(\theta) + \frac{1}{n} + C \left[\frac{r^2 b^2 \log n}{2\gamma^2} + \log \frac{n}{\delta} \right] \right) \end{aligned}$$

for all $\|\theta\|_2 \leq r$.

Rewriting (replacing 2γ with γ) and recognizing that with no loss of generality we may take γ such that $rb \geq \gamma$ gives the claim of the proposition. \square

6.2.3 A mutual information bound

An alternative perspective of the PAC-Bayesian bounds that Theorem 6.2.1 gives is to develop bounds based on mutual information, which is also central to the interactive data analysis setting in the next section. We present a few results along these lines here. Assume the setting of Theorem 6.2.1, so that \mathcal{F} consists of σ^2 -sub-Gaussian functions. Let us assume the following observational model: we observe $X_1^n \stackrel{\text{iid}}{\sim} P$, and then conditional on the sample X_1^n , draw a (random) function $F \in \mathcal{F}$ following the distribution $\pi(\cdot | X_1^n)$. Assuming the prior π_0 is fixed, Theorem 6.2.1 guarantees that with probability at least $1 - \delta$ over X_1^n ,

$$\mathbb{E}[(P_n F - P F)^2 | X_1^n] \leq \frac{8\sigma^2}{3n} \left[D_{\text{kl}}(\pi(\cdot | X_1^n) \| \pi_0) + \log \frac{2}{\delta} \right],$$

where the expectation is taken over $F \sim \pi(\cdot | X_1^n)$, leaving the sample fixed. Now, consider choosing π_0 to be the average over all samples X_1^n of π , that is, $\pi_0(\cdot) = \mathbb{E}_P[\pi(\cdot | X_1^n)]$, the expectation taken over $X_1^n \stackrel{\text{iid}}{\sim} P$. Then by definition of mutual information,

$$I(F; X_1^n) = \mathbb{E}_P [D_{\text{kl}}(\pi(\cdot | X_1^n) \| \pi_0)],$$

and by Markov’s inequality we have

$$\mathbb{P}(D_{\text{kl}}(\pi(\cdot | X_1^n) \| \pi_0) \geq K \cdot I(F; X_1^n)) \leq \frac{1}{K}$$

for all $K \geq 0$. Combining these, we obtain the following corollary.

Corollary 6.2.8. *Let F be chosen according to any distribution $\pi(\cdot | X_1^n)$ conditional on the sample X_1^n . Then with probability at least $1 - \delta_0 - \delta_1$ over the sample $X_1^n \stackrel{\text{iid}}{\sim} P$,*

$$\mathbb{E}[(P_n F - P F)^2 | X_1^n] \leq \frac{8\sigma^2}{3n} \left[\frac{I(F; X_1^n)}{\delta_0} + \log \frac{2}{\delta_1} \right].$$

This corollary shows that if we have any procedure—say, a learning procedure or otherwise—that limits the information between a sample X_1^n and an output F , then we are guaranteed that F generalizes. Tighter analyses of this are possible, though not our focus here, just that already there should be an inkling that limiting information between input samples and outputs may be fruitful.

6.3 Interactive data analysis

A major challenge in modern data analysis is that analyses are often not the classical statistics and scientific method setting. In the scientific method—forgive me for being a pedant—one proposes a hypothesis, the status quo or some other belief, and then designs an experiment to falsify that hypothesis. Then, upon performing the experiment, there are only two options: either the experimental results contradict the hypothesis (that is, we must reject the null) so that the hypothesis is false, or the hypothesis remains consistent with available data. In the classical (Fisherian) statistics perspective, this typically means that we have a single null hypothesis H_0 before observing a sample, we draw a sample $X \in \mathcal{X}$, and then for some test statistic $T : \mathcal{X} \rightarrow \mathbb{R}$ with observed value $t_{\text{observed}} = T(X)$, we compute the probability under the null of observing something as extreme as what we observed, that is, the p -value $p = P_{H_0}(T(X) \geq t_{\text{observed}})$.

Yet modern data analyses are distant from this pristine perspective for many reasons. The simplest is that we often have a number of hypotheses we wish to test, not a single one. For example, in biological applications, we may wish to investigate the associations between the expression of number of genes and a particular phenotype or disease; each gene j then corresponds to a null hypothesis $H_{0,j}$ that gene j is independent of the phenotype. There are numerous approaches to addressing the challenges associated with such multiple testing problems—such as false discovery rate control, familywise error rate control, and others—with whole courses devoted to the challenges.

Even these approaches to multiple testing and high-dimensional problems do not truly capture modern data analyses, however. Indeed, in many fields, researchers use one or a few main datasets, writing papers and performing multiple analyses on the same dataset. For example, in medicine, the UK Biobank dataset [174] has several thousand citations (as of 2023), many of which build on one another, with early studies coloring the analyses in subsequent studies. Even in situations without a shared dataset, analyses present researchers with huge degrees of freedom and choice. A researcher may study a summary statistic of his or her sampled data, or a plot of a few simple relationships, performing some simple data exploration—which statisticians and scientists have advocated for 50 years, dating back at least to John Tukey!—but this means that there are huge numbers of *potential* comparisons a researcher might make (that he or she does not). This “garden

of forking paths,” as Gelman and Loken [100] term it, causes challenges even when researchers are not “*p*-hacking” or going on a “fishing expedition” to try to find publishable results. The problem in these studies and approaches is that, because we make decisions that may, even only in a small way, depend on the data observed, we have invalidated all classical statistical analyses.

To that end, we now consider *interactive* data analyses, where we perform data analyses sequentially, computing new functions on a fixed sample X_1, \dots, X_n after observing some initial information about the sample. The starting point of our approach is similar to our analysis of PAC-Bayesian learning and generalization: we observe that if the function we decide to compute on the data X_1^n is chosen without much information about the data at hand, then its value on the sample should be similar to its values on the full population. This insight dovetails with what we have seen thus far, that appropriate “stability” in information can be useful and guarantee good future performance.

6.3.1 The interactive setting

We do not consider the interactive data analysis setting in full, rather, we consider a stylized approach to the problem, as it captures many of the challenges while being broad enough for different applications. In particular, we focus on the *statistical queries* setting, where a data analyst wishes to evaluate expectations

$$\mathbb{E}_P[\phi(X)] \tag{6.3.1}$$

of various functionals $\phi : \mathcal{X} \rightarrow \mathbb{R}$ under the population P using a sample $X_1^n \stackrel{\text{iid}}{\sim} P$. Certainly, numerous problems are solvable using statistical queries (6.3.1). Means use $\phi(x) = x$, while we can compute variances using the two statistical queries $\phi_1(x) = x$ and $\phi_2(x) = x^2$, as $\text{Var}(X) = \mathbb{E}_P[\phi_2(X)] - \mathbb{E}_P[\phi_1(X)]^2$.

Classical algorithms for the statistical query problem simply return sample means $P_n\phi := \frac{1}{n} \sum_{i=1}^n \phi(X_i)$ given a query $\phi : \mathcal{X} \rightarrow \mathbb{R}$. When the number of queries to be answered is not chosen adaptively, this means we can typically answer a large number relatively accurately; indeed, if we have a finite collection Φ of σ^2 -sub-Gaussian $\phi : \mathcal{X} \rightarrow \mathbb{R}$, then we of course have

$$\mathbb{P} \left(\max_{\phi \in \Phi} |P_n\phi - P\phi| \geq \sqrt{\frac{2\sigma^2}{n} (\log(2|\Phi|) + t)} \right) \leq e^{-t^2} \quad \text{for } t \geq 0$$

by Corollary 4.1.10 (sub-Gaussian concentration) and a union bound. Thus, so long as $|\Phi|$ is not exponential in the sample size n , we expect uniformly high accuracy.

Example 6.3.1 (Risk minimization via statistical queries): Suppose that we are in the loss-minimization setting (5.2.2), where the losses $\ell(\theta, X_i)$ are convex and differentiable in θ . Then gradient descent applied to $\hat{L}_n(\theta) = P_n\ell(\theta, X)$ will converge to a minimizing value of \hat{L}_n . We can evidently implement gradient descent by a sequence of statistical queries $\phi(x) = \nabla_{\theta}\ell(\theta, x)$, iterating

$$\theta^{(k+1)} = \theta^{(k)} - \alpha_k P_n\phi^{(k)}, \tag{6.3.2}$$

where $\phi^{(k)} = \nabla_{\theta}\ell(\theta^{(k)}, x)$ and α_k is a stepsize. \diamond

One issue with the example (6.3.1) is that we are *interacting* with the dataset, because each sequential query $\phi^{(k)}$ depends on the previous $k - 1$ queries. (Our results on uniform convergence of empirical functionals and related ideas address many of these challenges, so that the result of the process (6.3.2) will be well-behaved regardless of the interactivity.)

We consider an interactive version of the statistical query estimation problem. In this version, there are two parties: an analyst (or statistician or learner), who issues queries $\phi : \mathcal{X} \rightarrow \mathbb{R}$, and a mechanism that answers the queries to the analyst. We index our functionals ϕ by $t \in \mathcal{T}$ for a (possibly infinite) set \mathcal{T} , so we have a collection $\{\phi_t\}_{t \in \mathcal{T}}$. In this context, we thus have the following scheme:

Input: Sample X_1^n drawn i.i.d. P , collection $\{\phi_t\}_{t \in \mathcal{T}}$ of possible queries
Repeat: for $k = 1, 2, \dots$

- i. Analyst chooses index $T_k \in \mathcal{T}$ and query $\phi := \phi_{T_k}$
- ii. Mechanism responds with answer A_k approximating $P\phi = \mathbb{E}_P[\phi(X)]$ using X_1^n

Figure 6.1: The interactive statistical query setting

Of interest in the iteration 6.1 is that we *interactively* choose T_1, T_2, \dots, T_k , where the choice T_i may depend on our approximations of $\mathbb{E}_P[\phi_{T_j}(X)]$ for $j < i$, that is, on the results of our previous queries. Even more broadly, the analyst may be able to choose the index T_k in alternative ways depending on the sample X_1^n , and our goal is to still be able to accurately compute expectations $P\phi_T = \mathbb{E}_P[\phi_T(X)]$ when the index T may depend on X_1^n . The setting in Figure 6.1 clearly breaks with the classical statistical setting in which an analysis is pre-specified before collecting data, but more closely captures modern data exploration practices.

6.3.2 Second moment errors and mutual information

The starting point of our derivation is the following result, which follows from more or less identical arguments to those for our PAC-Bayesian bounds earlier.

Theorem 6.3.2. *Let $\{\phi_t\}_{t \in \mathcal{T}}$ be a collection of σ^2 -sub-Gaussian functions $\phi_t : \mathcal{X} \rightarrow \mathbb{R}$. Then for any random variable T and any $\lambda > 0$,*

$$\mathbb{E}[(P_n \phi_T - P \phi_T)^2] \leq \frac{1}{\lambda} \left[I(X_1^n; T) - \frac{1}{2} \log [1 - 2\lambda\sigma^2/n]_+ \right]$$

and

$$|\mathbb{E}[P_n \phi_T] - \mathbb{E}[P \phi_T]| \leq \sqrt{\frac{2\sigma^2}{n} I(X_1^n; T)}$$

where the expectations are taken over T and the sample X_1^n .

Proof The proof is similar to that of our first basic PAC-Bayes result in Theorem 6.2.1. Let us assume w.l.o.g. that $P\phi_t = 0$ for all $t \in \mathcal{T}$, noting that then $P_n \phi_t$ is σ^2/n -sub-Gaussian. We prove the first result first. Lemma 6.2.2 implies that $\mathbb{E}[\exp(\lambda(P_n \phi_t)^2)] \leq [1 - 2\lambda\sigma^2/n]_+^{-1/2}$ for each $t \in \mathcal{T}$. As a consequence, we obtain via the Donsker-Varadhan equality (Theorem 6.1.1) that

$$\begin{aligned} \lambda \mathbb{E} \left[\int (P_n \phi_t)^2 d\pi(t) \right] &\stackrel{(i)}{\leq} \mathbb{E}[D_{\text{kl}}(\pi \| \pi_0)] + \mathbb{E} \left[\log \int \exp(\lambda(P_n \phi_t)^2) d\pi_0(t) \right] \\ &\stackrel{(ii)}{\leq} \mathbb{E}[D_{\text{kl}}(\pi \| \pi_0)] + \log \mathbb{E} \left[\int \exp(\lambda(P_n \phi_t)^2) d\pi_0(t) \right] \\ &\stackrel{(iii)}{\leq} \mathbb{E}[D_{\text{kl}}(\pi \| \pi_0)] - \frac{1}{2} \log [1 - 2\lambda\sigma^2/n]_+ \end{aligned}$$

for all distributions π on \mathcal{T} , which may depend on P_n , where the expectation \mathbb{E} is taken over the sample $X_1^n \stackrel{\text{iid}}{\sim} P$. (Here inequality (i) is Theorem 6.1.1, inequality (ii) is Jensen's inequality, and inequality (iii) is Lemma 6.2.2.) Now, let π_0 be the marginal distribution on T (marginally over all observations X_1^n), and let π denote the posterior of T conditional on the sample X_1^n . Then $\mathbb{E}[D_{\text{kl}}(\pi \| \pi_0)] = I(X_1^n; T)$ by definition of the mutual information, giving the bound on the squared error.

For the second result, note that the Donsker-Varadhan equality implies

$$\lambda \mathbb{E} \left[\int P_n \phi_t d\pi(t) \right] \leq \mathbb{E}[D_{\text{kl}}(\pi \| \pi_0)] + \log \int \mathbb{E}[\exp(\lambda P_n \phi_t)] d\pi_0(t) \leq I(X_1^n; T) + \frac{\lambda^2 \sigma^2}{2n}.$$

Dividing both sides by λ gives $\mathbb{E}[P_n \phi_T] \leq \sqrt{2\sigma^2 I(X_1^n; T)/n}$, and performing the same analysis with $-\phi_T$ gives the second result of the theorem. \square

The key in the theorem is that if the mutual information—the Shannon information— $I(X; T)$ between the sample X and T is small, then the expected squared error can be small. To make this a bit clearer, let us choose values for λ in the theorem; taking $\lambda = \frac{n}{2e\sigma^2}$ gives the following corollary.

Corollary 6.3.3. *Let the conditions of Theorem 6.3.2 hold. Then*

$$\mathbb{E}[(P_n \phi_T - P \phi_T)^2] \leq \frac{2e\sigma^2}{n} I(X_1^n; T) + \frac{5\sigma^2}{4n}.$$

Consequently, if we can limit the amount of information any particular query T (i.e., ϕ_T) contains about the actual sample X_1^n , then guarantee reasonably high accuracy in the second moment errors $(P_n \phi_T - P \phi_T)^2$.

6.3.3 Limiting interaction in interactive analyses

Let us now return to the interactive data analysis setting of Figure 6.1, where we recall the stylized application of estimating mean functionals $P\phi$ for $\phi \in \{\phi_t\}_{t \in \mathcal{T}}$. To motivate a more careful approach, we consider a simple example to show the challenges that may arise even with only a single “round” of interactive data analysis. Naively answering queries accurately—using the mechanism $P_n \phi$ that simply computes the sample average—can easily lead to problems:

Example 6.3.4 (A stylized correlation analysis): Consider the following stylized genetics experiment. We observe vectors $X \in \{-1, 1\}^k$, where $X_j = 1$ if gene j is expressed and -1 otherwise. We also observe phenotypes $Y \in \{-1, 1\}$, where $Y = 1$ indicates appearance of the phenotype. In our setting, we will assume that the vectors X are uniform on $\{-1, 1\}^k$ and independent of Y , but an experimentalist friend of ours wishes to know if there exists a vector v with $\|v\|_2 = 1$ such that the correlation between $v^T X$ and Y is high, meaning that $v^T X$ is associated with Y . In our notation here, we have index set $\{v \in \mathbb{R}^k \mid \|v\|_2 = 1\}$, and by Example 4.1.6, Hoeffding's lemma, and the independence of the coordinates of X we have that $v^T X Y$ is $\|v\|_2^2/4 = 1/4$ -sub-Gaussian. Now, we recall the fact that if Z_j , $j = 1, \dots, k$, are σ^2 -sub-Gaussian, then for any $p \geq 1$, we have

$$\mathbb{E}[\max_j |Z_j|^p] \leq (Cp\sigma^2 \log k)^{p/2}$$

for a numerical constant C . That is, powers of sub-Gaussian maxima grow at most logarithmically. Indeed, by Theorem 4.1.11, we have for any $q \geq 1$ by Hölder's inequality that

$$\mathbb{E}[\max_j |Z_j|^p] \leq \mathbb{E}\left[\sum_j |Z_j|^{pq}\right]^{1/q} \leq k^{1/q} (Cpq\sigma^2)^{p/2},$$

and setting $q = \log k$ gives the inequality. Thus, we see that for any *a priori* fixed v_1, \dots, v_k, v_{k+1} , we have

$$\mathbb{E}[\max_j (v_j^T (P_n Y X))^2] \leq O(1) \frac{\log k}{n}.$$

If instead we allow a *single* interaction, the problem is different. We issue queries associated with $v = e_1, \dots, e_k$, the k standard basis vectors; then we simply set $V_{k+1} = P_n Y X / \|P_n Y X\|_2$. Then evidently

$$\mathbb{E}[(V_{k+1}^T (P_n Y X))^2] = \mathbb{E}[\|P_n Y X\|_2^2] = \frac{k}{n},$$

which is exponentially larger than in the non-interactive case. That is, if an analyst is allowed to interact with the dataset, he or she may be able to discover very large correlations that are certainly false in the population, which in this case has $PXY = 0$. \diamond

Example 6.3.4 shows that, without being a little careful, substantial issues may arise in interactive data analysis scenarios. When we consider our goal more broadly, which is to be able to provide accurate approximations to $P\phi$ for queries ϕ chosen adaptively for any population distribution P and $\phi : \mathcal{X} \rightarrow [-1, 1]$, it is possible to construct quite perverse situations, where if we compute sample expectations $P_n\phi$ exactly, one round of interaction is sufficient to find a query ϕ for which $P_n\phi - P\phi \geq 1$.

Example 6.3.5 (Exact query answering allows arbitrary corruption): Suppose we draw a sample X_1^n of size n on a sample space $\mathcal{X} = [m]$ with $X_i \stackrel{\text{iid}}{\sim} \text{Uniform}([m])$, where $m \geq 2n$. Let Φ be the collection of all functions $\phi : [m] \rightarrow [-1, 1]$, so that $\mathbb{P}(|P_n\phi - P\phi| \geq t) \leq \exp(-nt^2/2)$ for any fixed ϕ . Suppose that in the interactive scheme in Fig. 6.1, we simply release answers $A = P_n\phi$. Consider the following query:

$$\phi(x) = n^{-x} \text{ for } x = 1, 2, \dots, m.$$

Then by inspection, we see that

$$\begin{aligned} P_n\phi &= \sum_{j=1}^m n^{-j} \text{card}(\{X_i \mid X_i = j\}) \\ &= \frac{1}{n} \text{card}(\{X_i \mid X_i = 1\}) + \frac{1}{n^2} \text{card}(\{X_i \mid X_i = 1\}) + \dots + \frac{1}{n^m} \text{card}(\{X_i \mid X_i = m\}). \end{aligned}$$

It is clear that given $P_n\phi$, we can reconstruct the sample counts exactly. Then if we define a second query $\phi_2(x) = 1$ for $x \in X_1^n$ and $\phi_2(x) = -1$ for $x \notin X_1^n$, we see that $P\phi_2 \leq \frac{n}{m} - 1$, while $P_n\phi_2 = 1$. The gap is thus

$$\mathbb{E}[P_n\phi_2 - P\phi_2] \geq 2 - \frac{n}{m} \geq 1,$$

which is essentially as bad as possible. \diamond

More generally, when one performs an interactive data analysis (e.g. as in Fig. 6.1), adapting hypotheses while interacting with a dataset, it is not a question of statistical significance or multiplicity control for the analysis one does, but for *all the possible analyses* one might have done otherwise. Given the branching paths one might take in an analysis, it is clear that we require some care.

With that in mind, we consider the desiderata for techniques we might use to control information in the indices we select. We seek some type of *stability* in the information algorithms provide to a data analyst—intuitively, if small changes to a sample do not change the behavior of an analyst substantially, then we expect to obtain reasonable generalization bounds. If outputs of a particular analysis procedure carry little information about a particular sample (but instead provide information about a population), then Corollary 6.3.3 suggests that any estimates we obtain should be accurate.

To develop this stability theory, we require two conditions: first, that whatever quantity we develop for stability should *compose adaptively*, meaning that if we apply two (randomized) algorithms to a sample, then if both are appropriately stable, even if we choose the second algorithm because of the output of the first in arbitrary ways, they should remain jointly stable. Second, our notion should bound the mutual information $I(X_1^n; T)$ between the sample X_1^n and T . Lastly, we remark that this control on the mutual information has an additional benefit: by the data processing inequality, any downstream analysis we perform that depends only on T necessarily satisfies the same stability and information guarantees as T , because if we have the Markov chain $X_1^n \rightarrow T \rightarrow V$ then $I(X_1^n; V) \leq I(X_1^n; T)$.

We consider randomized algorithms $A : \mathcal{X}^n \rightarrow \mathcal{A}$, taking values in our index set \mathcal{A} , where $A(X_1^n) \in \mathcal{A}$ is a random variable that depends on the sample X_1^n . For simplicity in derivation, we abuse notation in this section, and for random variables X and Y with distributions P and Q respectively, we denote

$$D_{\text{kl}}(X \| Y) := D_{\text{kl}}(P \| Q).$$

We then ask for a type of leave-one-out stability for the algorithms A , where A is insensitive to the changes of a single example (on average).

Definition 6.1. Let $\varepsilon \geq 0$. A randomized algorithm $A : \mathcal{X}^n \rightarrow \mathcal{A}$ is ε -KL-stable if for each $i \in \{1, \dots, n\}$ there is a randomized $A_i : \mathcal{X}^{n-1} \rightarrow \mathcal{A}$ such that for every sample $x_1^n \in \mathcal{X}^n$,

$$\frac{1}{n} \sum_{i=1}^n D_{\text{kl}}(A(x_1^n) \| A_i(x_{\setminus i})) \leq \varepsilon.$$

Examples may be useful to understand Definition 6.1.

Example 6.3.6 (KL-stability in mean estimation: Gaussian noise addition): Suppose we wish to estimate a mean, and that $x_i \in [-1, 1]$ are all real-valued. Then a natural statistic is to simply compute $A(x_1^n) = \frac{1}{n} \sum_{i=1}^n x_i$. In this case, without randomization, we will have infinite KL-divergence between $A(x_1^n)$ and $A_i(x_{\setminus i})$. If instead we set $A(x_1^n) = \frac{1}{n} \sum_{i=1}^n x_i + Z$ for $Z \sim \mathcal{N}(0, \sigma^2)$, and similarly $A_i = \frac{1}{n} \sum_{j \neq i} x_j + Z$, then we have (recall Example 2.1.7)

$$\frac{1}{n} \sum_{i=1}^n D_{\text{kl}}(A(x_1^n) \| A_i(x_{\setminus i})) = \frac{1}{2n\sigma^2} \sum_{i=1}^n \frac{1}{n^2} x_i^2 \leq \frac{1}{2\sigma^2 n^2},$$

so that a the sample mean of a bounded random variable perturbed with Gaussian noise is $\varepsilon = \frac{1}{2\sigma^2 n^2}$ -KL-stable. \diamond

We can consider other types of noise addition as well.

Example 6.3.7 (KL-stability in mean estimation: Laplace noise addition): Let the conditions of Example 2.1.7 hold, but suppose instead of Gaussian noise we add scaled Laplace noise, that is, $A(x_1^n) = \frac{1}{n} \sum_{i=1}^n x_i + Z$ for Z with density $p(z) = \frac{1}{2\sigma} \exp(-|z|/\sigma)$, where $\sigma > 0$. Then using that if $L_{\mu,\sigma}$ denotes the Laplace distribution with shape σ and mean μ , with density $p(z) = \frac{1}{2\sigma} \exp(-|z - \mu|/\sigma)$, we have

$$\begin{aligned} D_{\text{kl}}(L_{\mu_0,\sigma} \| L_{\mu_1,\sigma}) &= \frac{1}{\sigma^2} \int_0^{|\mu_1 - \mu_0|} \exp(-z/\sigma) (|\mu_1 - \mu_0| - z) dz \\ &= \exp\left(-\frac{|\mu_1 - \mu_0|}{\sigma}\right) - 1 + \frac{|\mu_1 - \mu_0|}{\sigma} \leq \frac{|\mu_1 - \mu_0|^2}{2\sigma^2}, \end{aligned}$$

we see that in this case the sample mean of a bounded random variable perturbed with Laplace noise is $\varepsilon = \frac{1}{2\sigma^2 n^2}$ -KL-stable, where σ is the shape parameter. \diamond

The two key facts are that KL-stable algorithms compose adaptively and that they bound mutual information in independent samples.

Lemma 6.3.8. *Let $A : \mathcal{X}^n \rightarrow \mathcal{A}_0$ and $A' : \mathcal{A}_0 \times \mathcal{X} \rightarrow \mathcal{A}_1$ be ε and ε' -KL-stable algorithms, respectively. Then the (randomized) composition $A' \circ A(x_1^n) = A'(A(x_1^n), x_1^n)$ is $\varepsilon + \varepsilon'$ -KL-stable. Moreover, the pair $(A' \circ A(x_1^n), A(x_1^n))$ is $\varepsilon + \varepsilon'$ -KL-stable.*

Proof Let A_i and A'_i be the promised sub-algorithms in Definition 6.1. We apply the data processing inequality, which implies for each i that

$$D_{\text{kl}}(A'(A(x_1^n), x_1^n) \| A'_i(A_i(x_{\setminus i}), x_{\setminus i})) \leq D_{\text{kl}}(A'(A(x_1^n), x_1^n), A(x_1^n) \| A'_i(A_i(x_{\setminus i}), x_{\setminus i}), A_i(x_{\setminus i})).$$

We require a bit of notational trickery now. Fixing i , let $P_{A,A'}$ be the joint distribution of $A'(A(x_1^n), x_1^n)$ and $A(x_1^n)$ and $Q_{A,A'}$ the joint distribution of $A'_i(A_i(x_{\setminus i}), x_{\setminus i})$ and $A_i(x_{\setminus i})$, so that they are both distributions over $\mathcal{A}_1 \times \mathcal{A}_0$. Let $P_{A'|a}$ be the distribution of $A'(t, x_1^n)$ and similarly $Q_{A'|a}$ is the distribution of $A'_i(t, x_{\setminus i})$. Note that A', A'_i both “observe” x , so that using the chain rule (2.1.6) for KL-divergences, we have

$$\begin{aligned} D_{\text{kl}}(A' \circ A, A \| A'_i \circ A_i, A_i) &= D_{\text{kl}}(P_{A,A'} \| Q_{A,A'}) \\ &= D_{\text{kl}}(P_A \| Q_A) + \int D_{\text{kl}}(P_{A'|t} \| Q_{A'|t}) dP_A(t) \\ &= D_{\text{kl}}(A \| A_i) + \mathbb{E}_A[D_{\text{kl}}(A'(A, x_1^n) \| A'_i(A, x_1^n))]. \end{aligned}$$

Summing this from $i = 1$ to n yields

$$\frac{1}{n} \sum_{i=1}^n D_{\text{kl}}(A' \circ A \| A'_i \circ A_i) \leq \frac{1}{n} \sum_{i=1}^n D_{\text{kl}}(A \| A_i) + \mathbb{E}_A \left[\frac{1}{n} \sum_{i=1}^n D_{\text{kl}}(A'(A, x_1^n) \| A'_i(A, x_1^n)) \right] \leq \varepsilon + \varepsilon',$$

as desired. \square

The second key result is that KL-stable algorithms also bound the mutual information of a random function.

Lemma 6.3.9. *Let X_i be independent. Then for any random variable A ,*

$$I(A; X_1^n) \leq \sum_{i=1}^n I(A; X_i | X_{\setminus i}) = \sum_{i=1}^n \int D_{\text{kl}}(A(x_1^n) \| A_i(x_{\setminus i})) dP(x_1^n),$$

where $A_i(x_{\setminus i}) = A(x_1^{i-1}, X_i, x_{i+1}^n)$ is the random realization of A conditional on $X_{\setminus i} = x_{\setminus i}$.

Proof Without loss of generality, we assume A and X are both discrete. In this case, we have

$$I(A; X_1^n) = \sum_{i=1}^n I(A; X_i | X_1^{i-1}) = \sum_{i=1}^n H(X_i | X_1^{i-1}) - H(X_i | A, X_1^{i-1}).$$

Now, because the X_i follow a product distribution, $H(X_i | X_1^{i-1}) = H(X_i)$, while $H(X_i | A, X_1^{i-1}) \geq H(X_i | A, X_{\setminus i})$ because conditioning reduces entropy. Consequently, we have

$$I(A; X_1^n) \leq \sum_{i=1}^n H(X_i) - H(X_i | A, X_{\setminus i}) = \sum_{i=1}^n I(A; X_i | X_{\setminus i}).$$

To see the final equality, note that

$$\begin{aligned} I(A; X_i | X_{\setminus i}) &= \int_{\mathcal{X}^{n-1}} I(A; X_i | X_{\setminus i} = x_{\setminus i}) dP(x_{\setminus i}) \\ &= \int_{\mathcal{X}^{n-1}} \int_{\mathcal{X}} D_{\text{kl}}(A(x_1^n) \| A(x_{1:i-1}, X_i, x_{i+1:n})) dP(x_i) dP(x_{\setminus i}) \end{aligned}$$

by definition of mutual information as $I(X; Y) = \mathbb{E}_X[D_{\text{kl}}(P_{Y|X} \| P_Y)]$. □

Combining Lemmas 6.3.8 and 6.3.9, we see (nearly) immediately that KL stability implies a mutual information bound, and consequently even interactive KL-stable algorithms maintain bounds on mutual information.

Proposition 6.3.10. *Let A_1, \dots, A_k be ε_i -KL-stable procedures, respectively, composed in any arbitrary sequence. Let X_i be independent. Then*

$$\frac{1}{n} I(A_1, \dots, A_k; X_1^n) \leq \sum_{i=1}^k \varepsilon_i.$$

Proof Applying Lemma 6.3.9,

$$I(A_1^k; X_1^n) \leq \sum_{i=1}^n I(A_1^k; X_i | X_{\setminus i}) = \sum_{j=1}^k \sum_{i=1}^n I(A_j; X_i | X_{\setminus i}, A_1^{j-1}).$$

Fix an index j and for shorthand, let $A = A$ and $A' = (A_1, \dots, A_{j-1})$ be the first $j-1$ procedures. Then expanding the final mutual information term and letting ν denote the distribution of A' , we have

$$I(A; X_i | X_{\setminus i}, A') = \int D_{\text{kl}}(A(a', x_1^n) \| \bar{A}(a', x_{\setminus i})) dP(x_i | A' = a', x_{\setminus i}) dP^{n-1}(x_{\setminus i}) d\nu(a' | x_{\setminus i})$$

where $A(a', x_1^n)$ is the (random) procedure A on inputs x_1^n and a' , while $\bar{A}(a', x_{\setminus i})$ denotes the (random) procedure A on input $a', x_{\setminus i}, X_i$, and where the i th example X_i follows its distribution conditional on $A' = a'$ and $X_{\setminus i} = x_{\setminus i}$, as in Lemma 6.3.9. We then recognize that for each i , we have

$$\int D_{\text{kl}}(A(a', x_1^n) \| \bar{A}(a', x_{\setminus i})) dP(x_i | a', x_{\setminus i}) \leq \int D_{\text{kl}}(A(a', x_1^n) \| \tilde{A}(a', x_{\setminus i})) dP(x_i | a', x_{\setminus i})$$

for *any* randomized function \tilde{A} , as the marginal \bar{A} in the lemma minimizes the average KL-divergence (recall Exercise 2.15). Now, sum over i and apply the definition of KL-stability as in Lemma 6.3.8. \square

6.3.4 Error bounds for a simple noise addition scheme

Based on Proposition 6.3.10, to build an appropriately well-generalizing procedure we must build a mechanism for the interaction in Fig. 6.1 that maintains KL-stability. Using Example 6.3.6, this is not challenging for the class of bounded queries. Let $\Phi = \{\phi_t\}_{t \in \mathcal{T}}$ where $\phi_t : \mathcal{X} \rightarrow [-1, 1]$ be the collection of statistical queries taking values in $[-1, 1]$. Then based on Proposition 6.3.10 and Example 6.3.6, the following procedure is stable.

Input: Sample $X_1^n \in \mathcal{X}^n$ drawn i.i.d. P , collection $\{\phi_t\}_{t \in \mathcal{T}}$ of possible queries $\phi_t : \mathcal{X} \rightarrow [-1, 1]$

Repeat: for $k = 1, 2, \dots$

- i. Analyst chooses index $T_k \in \mathcal{T}$ and query $\phi := \phi_{T_k}$
- ii. Mechanism draws independent $Z_k \sim \mathcal{N}(0, \sigma^2)$ and responds with answer

$$A_k := P_n \phi + Z_k = \frac{1}{n} \sum_{i=1}^n \phi(X_i) + Z_k.$$

Figure 6.2: Sequential Gaussian noise mechanism.

This procedure is evidently KL-stable, and based on Example 6.3.6 and Proposition 6.3.10, we have that

$$\frac{1}{n} I(X_1^n; T_1, \dots, T_k, T_{k+1}) \leq \frac{k}{2\sigma^2 n^2}$$

so long as the indices $T_i \in \mathcal{T}$ are chosen only as functions of $P_n \phi + Z_j$ for $j < i$, as the classical information processing inequality implies that

$$\frac{1}{n} I(X_1^n; T_1, \dots, T_k, T_{k+1}) \leq \frac{1}{n} I(X_1^n; A_1, \dots, A_k)$$

because we have $X_1^n \rightarrow A_1 \rightarrow T_2$ and so on for the remaining indices. With this, we obtain the following theorem.

Theorem 6.3.11. *Let the indices T_i , $i = 1, \dots, k+1$ be chosen in an arbitrary way using the procedure 6.2, and let $\sigma^2 > 0$. Then*

$$\mathbb{E} \left[\max_{j \leq k} (A_j - P\phi_{T_j})^2 \right] \leq \frac{2ek}{\sigma^2 n^2} + \frac{10}{4n} + 4\sigma^2(\log k + 1).$$

By inspection, we can optimize over σ^2 by setting $\sigma^2 = \sqrt{k/(\log k + 1)}/n$, which yields the upper bound

$$\mathbb{E} \left[\max_{j \leq k} (A_j - P\phi_{T_j})^2 \right] \leq \frac{10}{4n} + 10 \frac{\sqrt{k(1 + \log k)}}{n}.$$

Comparing to Example 6.3.4, we see a substantial improvement. While we do not achieve accuracy scaling with $\log k$, as we would if the queried functionals ϕ_t were completely independent of the sample, we see that we achieve mean-squared error of order

$$\frac{\sqrt{k \log k}}{n}$$

for k adaptively chosen queries.

Proof To prove the result, we use a technique sometimes called the *monitor* technique. Roughly, the idea is that we can choose the index T_{k+1} in any way we desire as long as it is a function of the answers A_1, \dots, A_k and any other constants independent of the data. Thus, we may choose

$$T_{k+1} := T_{k^*} \text{ where } k^* = \operatorname{argmax}_{j \leq k} |A_j - P\phi_{T_j}|,$$

as this is a (downstream) function of the k different $\varepsilon = \frac{1}{2\sigma^2 n^2}$ -KL-stable queries T_1, \dots, T_k . As a consequence, we have from Corollary 6.3.3 (and the fact that the queries ϕ are 1-sub-Gaussian) that for $T = T_{k+1}$,

$$\mathbb{E}[(P_n \phi_T - P\phi_T)^2] \leq \frac{2e}{n} I(X_1^n; T_{k+1}) + \frac{5}{4n} \leq 2ek\varepsilon + \frac{5}{4n} = \frac{ek}{\sigma^2 n^2} + \frac{5}{4n}.$$

Now, we simply consider the independent noise addition, noting that $(a + b)^2 \leq 2a^2 + 2b^2$ for any $a, b \in \mathbb{R}$, so that

$$\begin{aligned} \mathbb{E} \left[\max_{j \leq k} (A_j - P\phi_{T_j})^2 \right] &\leq 2\mathbb{E}[(P_n \phi_T - P\phi_T)^2] + 2\mathbb{E} \left[\max_{j \leq k} \{Z_j^2\} \right] \\ &\leq \frac{2ek}{\sigma^2 n^2} + \frac{10}{4n} + 4\sigma^2(\log k + 1), \end{aligned} \tag{6.3.3}$$

where inequality (6.3.3) is the desired result and follows by the following lemma.

Lemma 6.3.12. *Let W_j , $j = 1, \dots, k$ be independent $\mathcal{N}(0, 1)$. Then $\mathbb{E}[\max_j W_j^2] \leq 2(\log k + 1)$.*

Proof We assume that $k \geq 3$, as the result is trivial otherwise. Using the tail bound for Gaussians (Mills's ratio for Gaussians, which is tighter than the standard sub-Gaussian bound) that $\mathbb{P}(W \geq t) \leq \frac{1}{\sqrt{2\pi}t} e^{-t^2/2}$ for $t \geq 0$ and that $\mathbb{E}[Z] = \int_0^\infty \mathbb{P}(Z \geq t) dt$ for a nonnegative random variable Z , we obtain that for any t_0 ,

$$\begin{aligned} \mathbb{E}[\max_j W_j^2] &= \int_0^\infty \mathbb{P}(\max_j W_j^2 \geq t) dt \leq t_0 + \int_{t_0}^\infty \mathbb{P}(\max_j W_j^2 \geq t) dt \\ &\leq t_0 + 2k \int_{t_0}^\infty \mathbb{P}(W_1 \geq \sqrt{t}) dt \leq t_0 + \frac{2k}{\sqrt{2\pi}} \int_{t_0}^\infty e^{-t/2} dt = t_0 + \frac{4k}{\sqrt{2\pi}} e^{-t_0/2}. \end{aligned}$$

Setting $t_0 = 2 \log(4k/\sqrt{2\pi})$ gives $\mathbb{E}[\max_j W_j^2] \leq 2 \log k + \log \frac{4}{\sqrt{2\pi}} + 1$. □

□

6.4 Bibliography and further reading

PAC-Bayes techniques originated with work of David McAllester [144, 145, 146], and we remark on his excellently readable tutorial [147]. The particular approaches we take to our proofs in Section 6.2 follow Catoni [48] and McAllester [146]. The PAC-Bayesian bounds we present, that simultaneously for *any* distribution π on \mathcal{F} , if $F \sim \pi$ then

$$\mathbb{E}[(P_n F - P F)^2 \mid X_1^n] \lesssim \frac{1}{n} \left[D_{\text{kl}}(\pi \parallel \pi_0) + \log \frac{1}{\delta} \right]$$

with probability at least $1 - \delta$ suggest that we can optimize them by choosing π carefully. For example, in the context of learning a statistical model parameterized by $\theta \in \Theta$ with losses $\ell(\theta; x, y)$, it is natural to attempt to find π minimizing

$$\mathbb{E}_\pi[P_n \ell(\theta; X, Y) \mid P_n] + C \sqrt{\frac{1}{n} D_{\text{kl}}(\pi \parallel \pi_0)}$$

in π , where the expectation is taken over $\theta \sim \pi$. If this quantity has optimal value ϵ_n^* , then one is immediately guaranteed that for the population P , we have $\mathbb{E}_\pi[P \ell(\theta; X, y)] \leq \epsilon_n^* + C \sqrt{\log \frac{1}{\delta}} / \sqrt{n}$. Langford and Caruana [131] take this approach, and Dziugaite and Roy [85] use it to give (the first) non-trivial bounds for deep learning models.

The questions of interactive data analysis begin at least several decades ago, perhaps most profoundly highlighted positively by Tukey's *Exploratory Data Analysis* [183]. Problems of scientific replicability have, conversely, highlighted many of the challenges of reusing data or peeking, even innocently, at samples before performing statistical analyses [118, 95, 100]. Our approach to formalizing these ideas, and making rigorous limiting information leakage, draws from a more recent strain of work in the theoretical computer science literature, with major contributions from Dwork, Feldman, Hardt, Pitassi, Reingold, and Roth and Bassily, Nissim, Smith, Steinke, Stemmer, and Ullman [84, 82, 83, 20, 21]. Our particular treatment most closely follows Feldman and Steinke [88]. The problems these techniques target also arise frequently in high-dimensional statistics, where one often wishes to estimate uncertainty and perform inference *after* selecting a model. While we do not touch on these problems, a few references in this direction include [27, 180, 114].

6.5 Exercises

Exercise 6.1 (Duality in Donsker-Varadhan): Here, we give a converse result to Theorem 6.1.1, showing that for any function $h : \mathcal{X} \rightarrow \mathbb{R}$,

$$\log \mathbb{E}_Q[e^{h(X)}] = \sup_P \{ \mathbb{E}_P[h(X)] - D_{\text{kl}}(P \parallel Q) \}, \quad (6.5.1)$$

where the supremum is taken over probability measures. If Q has a density, the supremum may be taken over probability measures having a density.

- (a) Show the equality (6.5.1) in the case that \mathcal{X} is discrete by directly computing the supremum. (That is, let $|\mathcal{X}| = k$, and identify probability measures P and Q with vectors $p, q \in \mathbb{R}_+^k$.)
- (b) Let Q have density q . Assume that $\mathbb{E}_Q[e^{h(X)}] < \infty$ and let

$$Z_h(x) = \exp(h(x)) / \mathbb{E}_Q[\exp(h(X))],$$

so $\mathbb{E}_Q[Z_h(X)] = 1$. Let P have density $p(x) = Z_h(x)q(x)$. Show that

$$\log \mathbb{E}_Q[e^{h(X)}] = \mathbb{E}_P[h(X)] - D_{\text{kl}}(P \| Q).$$

Why does this imply equality (6.5.1) in this case?

- (c) If $\mathbb{E}_Q[e^{h(X)}] = +\infty$, then monotone convergence implies that $\lim_{B \uparrow \infty} \mathbb{E}_Q[e^{\min\{B, h(X)\}}] = +\infty$. Conclude (6.5.1).

Exercise 6.2 (An alternative PAC-Bayes bound): Let $f : \Theta \times \mathcal{X} \rightarrow \mathbb{R}$, and let π_0 be a density on $\theta \in \Theta$. Use the dual form (6.5.1) of the variational representation of the KL-divergence show that with probability at least $1 - \delta$ over the draw of $X_1^n \stackrel{\text{iid}}{\sim} P$,

$$\int P_n f(\theta, X) \pi(\theta) d\theta \leq \int \log \mathbb{E}_P[\exp(f(\theta, X))] \pi(\theta) d\theta + \frac{D_{\text{kl}}(\pi \| \pi_0) + \log \frac{1}{\delta}}{n}$$

simultaneously for all distributions π on Θ , where the expectation \mathbb{E}_P is over $X \sim P$.

Exercise 6.3 (A mean estimator with sub-Gaussian concentration for a heavy-tailed distribution [49]): In this question, we use a PAC-Bayes bound to construct an estimator of the mean $\mathbb{E}[X]$ of a distribution with sub-Gaussian-like concentration that depends *only* on the second moments $\Sigma = \mathbb{E}[XX^\top]$ of the random vector X (not on any additional dimension-dependent quantities) while only assuming that $\mathbb{E}[\|X\|^2] < \infty$. Let ψ be an odd function (i.e., $\psi(-t) = -\psi(t)$) satisfying

$$-\log(1 - t + t^2) \leq \psi(t) \leq \log(1 + t + t^2).$$

The function $\psi(t) = \min\{1, \max\{-1, t\}\}$ (the truncation of t to the range $[-1, 1]$) is such a function. Let π_θ be the normal distribution $\mathcal{N}(\theta, \sigma^2 I)$ and π_0 be $\mathcal{N}(0, \sigma^2 I)$.

- (a) Let $\lambda > 0$. Use Exercise 6.2 to show that with probability at least $1 - \delta$, for all $\theta \in \mathbb{R}^d$

$$\frac{1}{\lambda} \int P_n \psi(\lambda \langle \theta', X \rangle) \pi_\theta(\theta') d\theta' \leq \langle \theta, \mathbb{E}[X] \rangle + \lambda \left(\theta^\top \Sigma \theta + \sigma^2 \text{tr}(\Sigma) \right) + \frac{\|\theta\|_2^2 / 2\sigma^2 + \log \frac{1}{\delta}}{n\lambda}.$$

- (b) For $\lambda > 0$, define the “directional mean” estimator

$$E_n(\theta, \lambda) = \frac{1}{\lambda} \int P_n \psi(\lambda \langle \theta', X \rangle) \pi_\theta(\theta') d\theta'.$$

Give a choice of $\lambda > 0$ such that with probability $1 - \delta$,

$$\sup_{\theta \in \mathbb{S}^{d-1}} |E_n(\theta, \lambda) - \langle \theta, \mathbb{E}[X] \rangle| \leq \frac{2}{\sqrt{n}} \sqrt{\left(\frac{1}{2\sigma^2} + \log \frac{1}{\delta} \right) \left(\|\Sigma\|_{\text{op}} + \sigma^2 \text{tr}(\Sigma) \right)},$$

where $\mathbb{S}^{d-1} = \{u \in \mathbb{R}^d \mid \|u\|_2 = 1\}$ is the unit sphere.

(c) Justify the following statement: choosing the vector $\hat{\mu}_n$ minimizing

$$\sup_{\theta \in \mathbb{S}^{d-1}} |E_n(\theta, \lambda) - \langle \theta, \mu \rangle|$$

in μ guarantees that with probability at least $1 - \delta$,

$$\|\hat{\mu}_n - \mathbb{E}[X]\|_2 \leq \frac{4}{\sqrt{n}} \sqrt{\left(\frac{1}{2\sigma^2} + \log \frac{1}{\delta}\right) (\|\Sigma\|_{\text{op}} + \sigma^2 \text{tr}(\Sigma))}.$$

(d) Give a choice of the prior/posterior variance σ^2 so that

$$\|\hat{\mu}_n - \mathbb{E}[X]\|_2 \leq \frac{4}{\sqrt{n}} \sqrt{\text{tr}(\Sigma) + 2 \|\Sigma\|_{\text{op}} \log \frac{1}{\delta}}$$

with probability at least $1 - \delta$.

Exercise 6.4 (Large-margin PAC-Bayes bounds for multiclass problems): Consider the following multiclass prediction scenario. Data comes in pairs $(x, y) \in b\mathbb{B}_2^d \times [k]$ where $\mathbb{B}_2^d = \{v \in \mathbb{R}^d \mid \|v\|_2 \leq 1\}$ denotes the ℓ_2 -ball and $[k] = \{1, \dots, k\}$. We make predictions using predictors $\theta_1, \dots, \theta_k \in \mathbb{R}^d$, where the prediction of y on an example x is

$$\hat{y}(x) := \underset{i \leq k}{\text{argmax}} \langle \theta_i, x \rangle.$$

We suffer an error whenever $\hat{y}(x) \neq y$, and the *margin* of our classifier on pair (x, y) is

$$\langle \theta_y, x \rangle - \max_{i \neq y} \langle \theta_i, x \rangle = \min_{i \neq y} \langle \theta_y - \theta_i, x \rangle.$$

If $\langle \theta_y, x \rangle > \langle \theta_i, x \rangle$ for all $i \neq y$, the margin is then positive (and the prediction is correct).

(a) Develop an analogue of the bounds in Section 6.2.2 in this k -class multiclass setting. To do so, you should (i) define the analogue of the margin-based loss ℓ_γ , (ii) show how Gaussian perturbations leave it similar, and (iii) prove an analogue of the bound in Section 6.2.2. You should assume one of the two conditions

$$(C1) \quad \|\theta_i\|_2 \leq r \text{ for all } i \quad (C2) \quad \sum_{i=1}^k \|\theta_i\|_2^2 \leq kr^2$$

on your classification vectors θ_i . Specify which condition you choose.

(b) Describe a minimization procedure—just a few lines suffice—that uses convex optimization to find a (reasonably) large-margin multiclass classifier.

Exercise 6.5 (A variance-based information bound): Let $\Phi = \{\phi_t\}_{t \in \mathcal{T}}$ be a collection of functions $\phi_t : \mathcal{X} \rightarrow \mathbb{R}$, where each ϕ_t satisfies the Bernstein condition (4.1.7) with parameters $\sigma^2(\phi_t)$ and b , that is, $|\mathbb{E}[(\phi_t(X) - P\phi_t(X))^k]| \leq \frac{k!}{2} \sigma^2(\phi_t) b^{k-2}$ for all $k \geq 3$ and $\text{Var}(\phi_t(X)) = \sigma^2(\phi_t)$. Let $T \in \mathcal{T}$ be any random variable, which may depend on an observed sample X_1^n . Show that for all $C > 0$ and $|\lambda| \leq \frac{C}{2b}$, then

$$\left| \mathbb{E} \left[\frac{P_n \phi_T - P \phi_T}{\max\{C, \sigma(\phi_T)\}} \right] \right| \leq \frac{1}{n|\lambda|} I(T; X_1^n) + |\lambda|.$$

Exercise 6.6 (An information bound on variance): Let $\Phi = \{\phi_t\}_{t \in \mathcal{T}}$ be a collection of functions $\phi_t : \mathcal{X} \rightarrow \mathbb{R}$, where each $\phi_t : \mathcal{X} \rightarrow [-1, 1]$. Let $\sigma^2(\phi_t) = \text{Var}(\phi_t(X))$. Let $s_n^2(\phi) = P_n \phi^2 - (P_n \phi)^2$ be the sample variance of ϕ . Show that for all $C > 0$ and $0 \leq \lambda \leq C/4$, then

$$\mathbb{E} \left[\frac{s_n^2(\phi_T)}{\max\{C, \sigma^2(\phi_T)\}} \right] \leq \frac{1}{n\lambda} I(T; X_1^n) + 2.$$

The $\max\{C, \sigma^2(\phi_T)\}$ term is there to help avoid division by 0. *Hint:* If $0 \leq x \leq 1$, then $e^x \leq 1 + 2x$, and if $X \in [0, 1]$, then $\mathbb{E}[e^X] \leq 1 + 2\mathbb{E}[X] \leq e^{2\mathbb{E}[X]}$. Use this to argue that $\mathbb{E}[e^{\lambda n P_n(\phi - P\phi)^2 / \max\{C, \sigma^2\}}] \leq e^{2\lambda n}$ for any $\phi : \mathcal{X} \rightarrow [-1, 1]$ with $\text{Var}(\phi) \leq \sigma^2$, then apply the Donsker-Varadhan theorem.

Exercise 6.7: Consider the following scenario: let $\phi : \mathcal{X} \rightarrow [-1, 1]$ and let $\alpha > 0$, $\tau > 0$. Let $\mu = P_n \phi$ and $s^2 = P_n \phi^2 - \mu^2$. Define $\sigma^2 = \max\{\alpha s^2, \tau^2\}$, and assume that $\tau^2 \geq \frac{5\alpha}{n}$.

(a) Show that the mechanism with answer A_k defined by

$$A := P_n \phi + Z \quad \text{for } Z \sim \mathcal{N}(0, \sigma^2)$$

is ε -KL-stable (Definition 6.1), where for a numerical constant $C < \infty$,

$$\varepsilon \leq C \cdot \frac{s^2}{n^2 \sigma^2} \cdot \left(1 + \frac{\alpha^2}{\sigma^2}\right).$$

(b) Show that if $\alpha^2 \leq C' \tau^2$ for a numerical constant $C' < \infty$, then we can take $\varepsilon \leq O(1) \frac{1}{n^2 \alpha}$.

Hint: Use exercise 2.14, and consider the “alternative” mechanisms of sampling from

$$\mathcal{N}(\mu_{-i}, \sigma_{-i}^2) \quad \text{where } \sigma_{-i}^2 = \max\{\alpha s_{-i}^2, \tau^2\}$$

for

$$\mu_{-i} = \frac{1}{n-1} \sum_{j \neq i} \phi(X_j) \quad \text{and} \quad s_{-i}^2 = \frac{1}{n-1} \sum_{j \neq i} \phi(X_j)^2 - \mu_{-i}^2.$$

Exercise 6.8 (A general variance-dependent bound on interactive queries): Consider the algorithm in Fig. 6.3. Let $\sigma^2(\phi_t) = \text{Var}(\phi_t(X))$ be the variance of ϕ_t .

(a) Show that for $b > 0$ and for all $0 \leq \lambda \leq \frac{b}{2}$,

$$\mathbb{E} \left[\max_{j \leq k} \frac{|A_j - P\phi_{T_j}|}{\max\{b, \sigma(\phi_{T_j})\}} \right] \leq \frac{1}{n\lambda} I(X_1^n; T_1^k) + \lambda + \sqrt{2 \log(ke)} \sqrt{\frac{4\alpha}{nb^2} I(X_1^n; T_1^k) + 2\alpha + \frac{\tau^2}{b^2}}.$$

(If you do not have quite the right constants, that's fine.)

(b) Using the result of Question 6.7, show that with appropriate choices for the parameters $\alpha, b, \tau^2, \lambda$ that for a numerical constant $C < \infty$

$$\mathbb{E} \left[\max_{j \leq k} \frac{|A_j - P\phi_{T_j}|}{\max\{(k \log k)^{1/4} / \sqrt{n}, \sigma(\phi_{T_j})\}} \right] \leq C \frac{(k \log k)^{1/4}}{\sqrt{n}}.$$

You may assume that k, n are large if necessary.

(c) Interpret the result from part (b). How does this improve over Theorem 6.3.11?

Input: Sample $X_1^n \in \mathcal{X}^n$ drawn i.i.d. P , collection $\{\phi_t\}_{t \in \mathcal{T}}$ of possible queries $\phi_t : \mathcal{X} \rightarrow [-1, 1]$, parameters $\alpha > 0$ and $\tau > 0$
Repeat: for $k = 1, 2, \dots$

- i. Analyst chooses index $T_k \in \mathcal{T}$ and query $\phi := \phi_{T_k}$
- ii. Set $s_k^2 := P_n \phi^2 - (P_n \phi)^2$ and $\sigma_k^2 := \max\{\alpha s_k^2, \tau^2\}$
- iii. Mechanism draws independent $Z_k \sim \mathcal{N}(0, \sigma_k^2)$ and responds with answer

$$A_k := P_n \phi + Z_k = \frac{1}{n} \sum_{i=1}^n \phi(X_i) + Z_k.$$

Figure 6.3: Sequential Gaussian noise mechanism with variance sensitivity.

Chapter 7

Advanced concentration inequalities

- I. Use Donsker-Varadhan to demonstrate basic equivalence of sub-Gaussianity and transport inequality [37, Corollar 4.14]
- II. Transport inequalities and their connections
- III. Potential idea: concentration without dimension for covariance matrices via PAC-Bayes bounds (i.e., matrix concentration from them)

Probably omit the entropy method stuff, focus instead on transportation

7.1 From divergences to concentration and back

The Donsker-Varadhan representation of the KL-divergence, Theorem 6.1.1, has a dual form, which allows us to provide an important connection between moment generating functions and the KL-divergence. We state this as a corollary, then connect it to moment generating function bounds.

Corollary 7.1.1. *Let P and Q be distributions on a common space \mathcal{X} . Then*

$$D_{\text{kl}}(P\|Q) = \sup_g \left\{ \mathbb{E}_P[g(X)] - \log \mathbb{E}_Q[e^{g(X)}] \right\},$$

where the supremum is taken over measurable functions $g : \mathcal{X} \rightarrow \mathbb{R}$ with $\mathbb{E}_Q[e^{g(X)}] < \infty$. Conversely, for any measurable $g : \mathcal{X} \rightarrow \mathbb{R}$ and distribution Q on \mathcal{X} ,

$$\log \mathbb{E}_Q[e^{g(X)}] = \sup_P \{ \mathbb{E}_P[g(X)] - D_{\text{kl}}(P\|Q) \},$$

where the supremum is taken over probability distributions P on \mathcal{X} with $\mathbb{E}_P[g(X)] < \infty$.

Proof The first claim is simply Theorem 6.1.1. For the second, we assume that $\mathbb{E}_Q[e^{g(X)}] < \infty$. (See Exercise 7.1 for the case that $\mathbb{E}_Q[e^{g(X)}] = +\infty$.) Then via the first part of the corollary, for any distribution P for which $\mathbb{E}_P[g(X)] < \infty$ we have $D_{\text{kl}}(P\|Q) \geq \mathbb{E}_P[g(X)] - \log \mathbb{E}_Q[e^{g(X)}]$, that is,

$$\log \mathbb{E}_Q[e^{g(X)}] \geq \mathbb{E}_P[g(X)] - D_{\text{kl}}(P\|Q).$$

To obtain equality, assume w.l.o.g. that Q has density q , and define P to have density $p(x) = \frac{e^{g(x)}}{\mathbb{E}_Q[e^{g(X)}]} q(x)$. This evidently integrates to 1, and

$$\begin{aligned} \mathbb{E}_P[g(X)] - D_{\text{kl}}(P\|Q) &= \frac{\mathbb{E}_Q[g(X)e^{g(X)}]}{\mathbb{E}_Q[e^{g(X)}]} - \mathbb{E}_Q \left[\frac{e^{g(X)}}{\mathbb{E}_Q[e^{g(X)}]} \log \frac{e^{g(X)}}{\mathbb{E}_Q[e^{g(X)}]} \right] \\ &= \log \mathbb{E}_Q[e^{g(X)}], \end{aligned}$$

completing the proof. \square

The key consequence of Corollary 7.1.1, for our purposes, is that it yields the first building block for our development of *transportation inequalities*, which relate concentration and moment generating functions to KL-divergence measures. The first of these uses the variational representation to provide an alternative characterization of sub-Gaussian random variables.

Theorem 7.1.2. *Let X be a real-valued random variable. Then the following are equivalent:*

(i) For all $\lambda \geq 0$,

$$\log \mathbb{E}_P[e^{\lambda(X - \mathbb{E}_P[X])}] \leq \frac{\lambda^2 \sigma^2}{2}.$$

(ii) For any probability distribution Q for which $D_{\text{kl}}(Q\|P) < \infty$,

$$\mathbb{E}_Q[X] - \mathbb{E}_P[X] \leq \sqrt{2\sigma^2 D_{\text{kl}}(Q\|P)}.$$

Proof To see that (i) implies (ii), note that swapping the roles of P and Q in Corollary 7.1.1 and taking $g(X) = \lambda(X - \mathbb{E}_P[X])$ yields that

$$\lambda(\mathbb{E}_Q[X] - \mathbb{E}_P[X]) - D_{\text{kl}}(Q\|P) \leq \frac{\lambda^2 \sigma^2}{2}$$

for all distributions Q and all $\lambda \geq 0$. Rearranging, we obtain

$$\mathbb{E}_Q[X] - \mathbb{E}_P[X] \leq \frac{D_{\text{kl}}(Q\|P)}{\lambda} + \frac{\lambda \sigma^2}{2},$$

and optimizing over $\lambda > 0$ yields (ii).

For the opposite direction, because $\inf_{\lambda > 0} \frac{a}{\lambda} + b\lambda = 2\sqrt{ab}$ for $a, b \geq 0$, we see that (ii) implies

$$\mathbb{E}_Q[X] - \mathbb{E}_P[X] \leq \frac{D_{\text{kl}}(Q\|P)}{\lambda} + \frac{\lambda \sigma^2}{2}$$

for all $\lambda > 0$. Rewriting this by taking $g(X) = X - \mathbb{E}_P[X]$, we have $\lambda \mathbb{E}_Q[g(X)] - D_{\text{kl}}(Q\|P) \leq \frac{\lambda^2 \sigma^2}{2}$ for all $\lambda > 0$ and Q . Applying Corollary 7.1.1 and taking a supremum over Q yields (i). \square

Theorem 7.1.2 provides another proof that bounded random variables are sub-Gaussian; the sheer slickness of the argument, which relies only on Pinsker's inequality, hints at the power of the approaches we will be able to develop. (Compare this approach to Example 4.1.6 and Exercise 4.1.)

Corollary 7.1.3. *Let X be a random variable taking values in $[a, b]$. Then X is $\frac{(b-a)^2}{4}$ -sub-Gaussian.*

Proof Assume P and Q have densities p and q w.r.t. a base measure μ . Then

$$\begin{aligned}\mathbb{E}_Q[X] - \mathbb{E}_P[X] &= \int x(q(x) - p(x))d\mu(x) \leq b \int_{q>p} (q(x) - p(x))d\mu(x) + a \int_{p>q} (q(x) - p(x))d\mu(x) \\ &= (b - a) \|P - Q\|_{\text{TV}},\end{aligned}$$

where we have used one of our many characterizations of the total variation (Lemma 2.2.4). Applying Pinsker's inequality (Proposition 2.2.8), we obtain

$$\mathbb{E}_Q[X] - \mathbb{E}_P[X] \leq (b - a) \sqrt{\frac{1}{2} D_{\text{kl}}(Q\|P)} = \sqrt{\frac{(b - a)^2}{2} D_{\text{kl}}(Q\|P)}.$$

This implies that f is $\frac{(b-a)^2}{4}$ -sub-Gaussian. \square

Pinsker's inequality that $\|P - Q\|_{\text{TV}}^2 \leq \frac{1}{2} D_{\text{kl}}(Q\|P)$ is the first example of a general family of *transportation inequalities*, which relate various distances on probability measures to the KL-divergence. We develop these inequalities and a few representative consequences in Section 7.2. Before doing so, however, we present an application of the variational representation to estimating a covariance as well as presenting a generalization of Theorem 7.1.2 beyond sub-Gaussian random variables.

7.1.1 Concentration of covariance matrices via the variational representation

The variational representation of the KL-divergence in Corollary 7.1.1 also allows proofs of strong concentration guarantees for random matrices, a subject of its own considerable interest. Typical approaches to such concentration guarantees involve the matrix Chernoff bound approaches we outline in Chapter 4.3, or approaches via covering numbers, as we present in Proposition 5.1.11. In this section, we provide an alternative approach that avoids these techniques, attacking the covariance more directly. To present the bounds, we use the following consequence of the variational representation of the KL-divergence, which follows from Corollary 7.1.1.

Corollary 7.1.4. *Let $f : \Theta \times \mathcal{X} \rightarrow \mathbb{R}$ and π_0 be a density on $\theta \in \Theta$. Then for any $t \geq 0$, with probability at least $1 - e^{-t}$ over the draw $X_1^n \stackrel{\text{iid}}{\sim} P$,*

$$\int P_n f(\theta, X) \pi(\theta) d\theta \leq \int \log \mathbb{E}_P [\exp(f(\theta, X))] \pi(\theta) d\theta + \frac{D_{\text{kl}}(\pi\|\pi_0) + t}{n}$$

simultaneously for all distributions π on Θ , where the expectation \mathbb{E}_P is over $X \sim P$.

Exercise 6.2 asks you to prove this result, which follows by an application of Markov's inequality to the nonnegative random variable $\exp(nP_n f(\theta, X)) = \prod_{i=1}^n \exp(f(\theta, X_i))$.

Returning to the study of random vectors and covariance matrices, here we provide an alternative argument to show that if $X_i \in \mathbb{R}^d$ are well-behaved random vectors, then the covariance or second-moment matrix

$$P_n X X^\top = \frac{1}{n} \sum_{i=1}^n X_i X_i^\top$$

concentrates well. We focus on d -dimensional σ^2 -sub-Gaussian vectors meaning that $\mathbb{E}[e^{\langle u, X - \mathbb{E}[X] \rangle}] \leq \exp(\frac{\sigma^2}{2} \|u\|_2^2)$ for all $u \in \mathbb{R}^d$. It will be more convenient to use the Orlicz-norm-based characterization of sub-Gaussianity (recall Section 4.1.3), so we say X is a σ^2 -sub-Gaussian vector if

$$\|\langle X, u \rangle\|_{\psi_2} \leq \sigma^2$$

for all $u \in \mathbb{S}^{d-1} = \{u \in \mathbb{R}^d \mid \|u\|_2 = 1\}$. (This is more convenient for analyzing $P_n X X^\top$, because the sub-multiplicity of the Orlicz norms gives $\|v^\top X X^\top u\|_{\psi_1} \leq \|\langle X, u \rangle\|_{\psi_2} \|\langle X, v \rangle\|_{\psi_2}$ as in Lemma 4.1.23.)

We then have the following proposition, which provides the same guarantees as Proposition 5.1.11 for isotropic X_i , meaning $\mathbb{E}[X_i X_i^\top] = I_d$.

Proposition 7.1.5. *Let X_i be independent isotropic sub-Gaussian vectors with $\|\langle u, X_i \rangle\|_{\psi_2} \leq \sigma$ for all $u \in \mathbb{S}^{d-1}$. Then with probability at least $1 - e^{-t}$,*

$$\|P_n X X^\top - I_d\|_{\text{op}} \leq \frac{8\epsilon^2 \sigma^2}{(\epsilon - 1)^2} \sqrt{\frac{2d + t}{n}}$$

so long as $n \geq 2d + t$. Otherwise, $\|P_n X X^\top - I_d\|_{\text{op}} \leq \frac{4\epsilon^2 \sigma^2}{(\epsilon - 1)^2} (1 + \frac{2d+t}{n})$.

Proof Fix $\epsilon > 0$ to be chosen, and let the prior density π_0 be uniform on the Cartesian product $(1 + \epsilon)\mathbb{B}_2^d \times (1 + \epsilon)\mathbb{B}_2^d$. For any fixed u, v with $\|u\|_2, \|v\|_2 \leq 1$, let the posterior $\pi_{u,v}$ be uniform on the product $(u + \epsilon\mathbb{B}_2^d) \times (v + \epsilon\mathbb{B}_2^d)$, so that for any vector x we have

$$\int \theta_1^\top x x^\top \theta_2 \pi_{u,v}(\theta_1, \theta_2) d\theta_1 d\theta_2 = u^\top x x^\top v.$$

Fix $\lambda \geq 0$ to be chosen, and define the function

$$f(\theta_1, \theta_2, x) := \theta_1^\top x x^\top \theta_2 - \langle \theta_1, \theta_2 \rangle.$$

Notably, f is mean zero under the distribution P , as $\mathbb{E}[X X^\top] = I_d$. Because we have assumed $\|\langle u, X \rangle\|_{\psi_2} \leq \sigma^2$ for all $u \in \mathbb{S}^{d-1}$, for $\theta_1, \theta_2 \in (1 + \epsilon)\mathbb{B}_2^d$ $f(\theta_1, \theta_2, X)$ is sub-exponential, and more precisely,

$$\begin{aligned} \|f(\theta_1, \theta_2, X)\|_{\psi_1} &= \left\| \theta_1^\top X X^\top \theta_2 - \langle \theta_1, \theta_2 \rangle \right\|_{\psi_1} \leq 2 \left\| \theta_1^\top X X^\top \theta_2 \right\|_{\psi_1} \\ &\leq 2 \|\langle \theta_1, X \rangle\|_{\psi_2} \|\langle \theta_2, X \rangle\|_{\psi_2} \leq 2(1 + \epsilon)^2 \sigma^2 \end{aligned} \quad (7.1.1)$$

by the triangle inequality and sub-multiplicativity (Lemma 4.1.23).

Let $\lambda \in \mathbb{R}$ and apply Corollary 7.1.4 to the function λf , which implies

$$\lambda \left(u^\top P_n X X^\top v - \langle u, v \rangle \right) \leq \frac{D_{\text{kl}}(\pi_{u,v} \| \pi_0) + t}{n} + \log \int \mathbb{E}_P[\exp(\lambda f(\theta_1, \theta_2, X))] \pi_0(\theta_1, \theta_2) d\theta_1 d\theta_2$$

simultaneously for all $u, v \in \mathbb{B}_2^d$ with probability at least $1 - e^{-t}$, where we recall that π is uniform on $(u + \epsilon\mathbb{B}_2^d) \times (v + \epsilon\mathbb{B}_2^d)$. It remains to control the expectation $\mathbb{E}_P[\exp(\lambda f)]$ and the KL-divergence. For the former, we apply Corollary 4.1.24 and use inequality (7.1.1) to guarantee that f has finite ψ_1 -norm, so for fixed $\theta_1, \theta_2 \in (1 + \epsilon)\mathbb{B}_2^d$ we have

$$\mathbb{E}_P[\exp(\lambda f(\theta_1, \theta_2, X))] \leq \exp(16\lambda^2(1 + \epsilon)^4 \sigma^4)$$

when $|\lambda| \leq \frac{1}{4(1+\epsilon)^2\sigma^2}$. Substituting this above

$$\lambda \left(u^\top P_n X X^\top v - \langle u, v \rangle \right) \leq \frac{D_{\text{kl}}(\pi \| \pi_0) + t}{n} + 16\lambda^2(1+\epsilon)^4\sigma^4.$$

Finally, we evaluate the KL-divergence: for any $u, v \in \mathbb{B}_2^d$, we have

$$D_{\text{kl}}(\pi_{u,v} \| \pi_0) = 2D_{\text{kl}}\left(\text{Uniform}(u + \epsilon \mathbb{B}_2^d) \| \text{Uniform}((1+\epsilon)\mathbb{B}_2^d)\right) = 2d \log \frac{1+\epsilon}{\epsilon}$$

because $\pi_{u,v}$ and π_0 are product distributions. Take $\epsilon = \frac{1}{e-1}$ to give $D_{\text{kl}}(\pi \| \pi_0) = 2d$ and $1+\epsilon = \frac{e}{e-1}$. Dividing by λ (accounting for the sign as appropriate), we see that with probability at least $1 - e^{-t}$, simultaneously for all $\|u\|_2 = 1$ and $\|v\|_2 = 1$ we have

$$|u^\top P_n X X^\top v - \langle u, v \rangle| \leq \frac{2d+t}{n\lambda} + 16 \left(\frac{e}{e-1} \right)^4 \sigma^4 \lambda$$

whenever $0 \leq \lambda \leq \frac{(e-1)^2}{4e^2\sigma^2}$. Take $\lambda = \frac{(e-1)^2}{4e^2\sigma^2} \min\{\sqrt{\frac{2d+t}{n}}, 1\}$. \square

Proposition 7.1.5 admits elegant extensions to non-isotropic matrices, giving “dimension-free” concentration for sub-Gaussian random vectors. Exercise 7.4 explores one approach to this result.

7.1.2 A generalized connection between moment generating functions and divergence

The approach we use to prove Theorem 7.1.2 extends beyond sub-Gaussian random variables to, essentially, any random variable whose moment generating function is even a bit nice, using (essentially) the same argument, because any log moment generating function has very similar variational properties to the quadratic (because, essentially, they look locally quadratic). Here, we elucidate this approach. Recall that for random variable X with log moment generating function $\phi_X(\lambda) := \log \mathbb{E}[e^{\lambda X}]$ (the *cumulant generating function*), so long as $\phi_X(\lambda)$ is finite in some neighborhood of 0, then it is infinitely differentiable on the interior of its domain (Proposition 3.2.2), and we may take derivatives through expectations. Thus, we always have

$$\phi'_X(\lambda) = \frac{\mathbb{E}[X e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} \quad \text{and} \quad \phi''_X(\lambda) = \frac{\mathbb{E}[X^2 e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} - \left(\frac{\mathbb{E}[X e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} \right)^2 = \text{Var}_{P_\lambda}(X),$$

where P_λ denotes the distribution on X with density tilted by $\frac{e^{\lambda x}}{\mathbb{E}[e^{\lambda X}]}$. So long as X is mean zero and non-constant, then, we obtain $\phi'_X(0) = 0$ and $\phi''_X(\lambda) > 0$ for all $\lambda \in \text{int dom } \phi_X$. (We can say a bit more about such cumulant generating functions; see Exercise 7.3.)

To that end, we call a function ϕ *CGF* (*cumulant-generating-function*) *like* on the interval $[0, b)$ if ϕ is continuously differentiable on $[0, b)$, satisfies $\phi(0) = \phi'(0) = 0$, and is convex (here we take $\phi'(0)$ to be the right derivative). The *convex conjugate* of ϕ is

$$\phi^*(s) := \sup_{0 \leq \lambda < b} \{\lambda s - \phi(\lambda)\},$$

which is strictly increasing in s . (We will have much more to say about convex conjugates in the coming chapters; see also Appendix B, especially Section C.2.) For now, we only require the generalized inverse

$$(\phi^*)^{-1}(t) := \inf \{s \geq 0 \mid \phi^*(s) > t\}.$$

Because $\phi^*(s) \geq \lambda s - \phi(\lambda)$ for some $\lambda > 0$, we have $\phi^*(s) \rightarrow \infty$ as $s \rightarrow \infty$, meaning that $(\phi^*)^{-1}(t)$ exists and is finite.

Example 7.1.6 (Sub-Gaussian CGF-like functions): If X is a σ^2 -sub-Gaussian random variables, we have $\phi_X(\lambda) \leq \frac{\lambda^2 \sigma^2}{2}$. We thus consider functions $\phi(\lambda) = \frac{\lambda^2 \sigma^2}{2}$, where $b = +\infty$. Then for $s \geq 0$, we have $\phi^*(s) = \sup\{\lambda s - \frac{\lambda^2 \sigma^2}{2}\} = \frac{s^2}{2\sigma^2}$ for $s \geq 0$, and

$$(\phi^*)^{-1}(t) = \sqrt{2\sigma^2 t} = \inf_{\lambda > 0} \frac{\lambda \sigma^2}{2} + t\lambda$$

gives the inverse. (Recall the proof of Theorem 7.1.2, which relies on this transformation.) \diamond

Example 7.1.7 (Sub-exponential CGF-like functions): When we have (τ^2, b) -sub-Exponential random variables (Definition 4.2), we obtain CGF-like functions of the form $\phi(\lambda) = \frac{\tau^2 \lambda^2}{2}$ for $0 \leq \lambda < \frac{1}{b}$. Taking suprema gives the the convex conjugate

$$\phi^*(s) = \begin{cases} \frac{s^2}{2\tau^2} & \text{if } s \leq \frac{\tau^2}{b} \\ \frac{s}{b} - \frac{\tau^2}{2b^2} & \text{if } s \geq \frac{\tau^2}{b}. \end{cases}$$

In this case, we have the more complex inverse $(\phi^*)^{-1}(t) = \min\{\sqrt{2\tau^2 t}, bt + \frac{\tau^2}{2b}\}$. \diamond

We can now present the generalization of Theorem 7.1.2.

Theorem 7.1.8. *Let X be a real-valued random variable and ϕ be CGF-like on $[0, b)$. Then the following are equivalent:*

(i) *For all $\lambda \in (0, b)$,*

$$\log \mathbb{E}_P[e^{\lambda(X - \mathbb{E}_P[X])}] \leq \phi(\lambda).$$

(ii) *For any probability distribution Q for which $D_{\text{kl}}(Q \| P) < \infty$,*

$$\mathbb{E}_Q[X] - \mathbb{E}_P[X] \leq (\phi^*)^{-1}(D_{\text{kl}}(Q \| P)).$$

Proof To see that (i) implies (ii), note that if $\phi(\lambda) \geq \log \mathbb{E}_P[e^{\lambda(X - \mathbb{E}_P[X])}]$, then as in the proof of Theorem 7.1.2, swapping the roles of P and Q in Corollary 7.1.1 and taking $g(X) = X - \mathbb{E}_P[X]$ yields that

$$\lambda(\mathbb{E}_Q[X] - \mathbb{E}_P[X]) - \phi(\lambda) \leq D_{\text{kl}}(Q \| P)$$

for all distributions Q . Taking a supremum over $\lambda \in (0, b)$ on the left side then gives

$$\phi^*(\mathbb{E}_Q[X] - \mathbb{E}_P[X]) \leq D_{\text{kl}}(Q \| P),$$

which implies part (ii).

For the converse direction, we require a technical lemma giving an inverse for convex conjugates of smooth increasing functions:

Lemma 7.1.9. *Let ϕ be CGF-like on $[0, b)$. Then*

$$(\phi^*)^{-1}(t) := \inf\{s \geq 0 \mid \phi^*(s) > t\} = \inf_{\lambda \in (0, b)} \frac{t + \phi(\lambda)}{\lambda}.$$

Additionally, $(\phi^)^{-1}$ is concave and strictly increasing on \mathbb{R}_+ , with $(\phi^*)^{-1}(t) \rightarrow \infty$ as $t \rightarrow \infty$.*

Proof We prove the equality first. As in our discussion earlier, because $\phi^*(s) \geq \lambda s - \phi(\lambda)$ for any $\lambda \in (0, b)$, the set $\{s \geq 0 \mid \phi^*(s) > t\}$ is non-empty. Then $\phi^*(s) > t$ if and only if

$$\lambda s - \phi(\lambda) > t$$

for some $\lambda \in (0, b)$, that is, $s > \frac{t + \phi(\lambda)}{\lambda}$ for some $\lambda \in (0, b)$.

The concavity follows because $(\phi^*)^{-1}(t)$ is the infimum of linear functions of t . Because $\inf_{\lambda > 0} \phi(\lambda)/\lambda = \phi'(0) = 0$ by convexity (the slope $\phi'(\lambda)$ is non-decreasing in λ), we have $(\phi^*)^{-1}(t) \geq \inf_{\lambda \in (0, b)} t/\lambda = t/b$ for all $t \geq 0$. Coupled with concavity this implies $(\phi^*)^{-1}$ is strictly increasing, with $\lim_{t \rightarrow \infty} (\phi^*)^{-1}(t) = \infty$. \square

Returning to the main thread of the argument, item (ii) then evidently implies

$$\mathbb{E}_Q[X] - \mathbb{E}_P[X] \leq \frac{D_{\text{kl}}(Q\|P) + \phi(\lambda)}{\lambda}$$

for all $\lambda \in (0, b)$, or, taking $g(X) = X - \mathbb{E}_P[X]$, that $\lambda \mathbb{E}_Q[g(X)] - D_{\text{kl}}(Q\|P) \leq \phi(\lambda)$ for all Q . Applying Corollary 7.1.1 and taking a supremum over Q then yields $\log \mathbb{E}_P[e^{\lambda g(X)}] \leq \phi(\lambda)$, as desired. \square

7.2 Transportation inequalities

In Corollary 7.1.3, we saw our first example of a *transportation inequality*: because

$$\|P - Q\|_{\text{TV}} \leq \sqrt{\frac{1}{2} D_{\text{kl}}(Q\|P)}$$

for any distributions P and Q , we saw that for any function f taking values in $[0, 1]$ we had

$$\mathbb{E}_Q[f] - \mathbb{E}_P[f] \leq \|P - Q\|_{\text{TV}} \leq \sqrt{\frac{1}{2} D_{\text{kl}}(Q\|P)},$$

and so *a fortiori* f is sub-Gaussian. We also have the characterization that $\sup_{f \in [0, 1]} \mathbb{E}_Q[f] - \mathbb{E}_P[f] = \|P - Q\|_{\text{TV}}$. To see this as a “transportation” inequality requires a bit more work.

JCD Comment: Perhaps better to just directly do Lagrangian calculation here.

By inspection, for *any* joint distribution π on X and Y with the correct marginals P and Q and any function f taking values in $[0, 1]$, we have $f(x) - f(y) \leq \mathbf{1}\{x \neq y\}$ and

$$\pi(X \neq Y) = \mathbb{E}_\pi[\mathbf{1}\{X \neq Y\}] \geq \mathbb{E}_\pi[f(X) - f(Y)] = \mathbb{E}_P[f(X)] - \mathbb{E}_Q[f(Y)].$$

Consider the special case that P and Q are Bernoulli distributions with parameters $p > q > 0$, so that $\|P - Q\|_{\text{TV}} = p - q$. Now, for $X \sim P$ and $Y \sim Q$, consider the joint distribution π for which

$$\pi := \begin{cases} X = 1, Y = 1 & \text{w.p. } q \\ X = 1, Y = 0 & \text{w.p. } p - q \\ X = 0, Y = 0 & \text{w.p. } 1 - p. \end{cases}$$

Then we evidently have $\pi(X \neq Y) = \pi(X = 1, Y = 0) = p - q = \|P - Q\|_{\text{TV}}$. In this case, at least, taking infima over joint distributions maintaining the marginals on X and Y ,

$$\inf_{\pi} \mathbb{E}_{\pi}[\mathbf{1}\{X \neq Y\}] = \|P - Q\|_{\text{TV}} = \sup_{f \in [0,1]} \{\mathbb{E}_P[f(X)] - \mathbb{E}_Q[f(Y)]\},$$

and we can transform this into a concentration inequality.

The preceding equality ends up holding in substantially more generality and yielding a plethora of powerful concentration inequalities. Given distributions P and Q on sets \mathcal{X} and \mathcal{Y} , we let $\Pi(P, Q)$ denote the set of *couplings* between the distributions, meaning joint distributions on $\mathcal{X} \times \mathcal{Y}$ whose marginals are correct:

$$\Pi(P, Q) := \{\pi \text{ on } \mathcal{X} \times \mathcal{Y} \text{ s.t. } \pi(\cdot, \mathcal{Y}) = P(\cdot) \text{ and } \pi(\mathcal{X}, \cdot) = Q(\cdot)\}.$$

Then for a nonnegative cost function $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, the *transportation cost* between P and Q is

$$W_c(P, Q) := \inf_{\pi \in \Pi(P, Q)} \mathbb{E}_{\pi}[c(X, Y)]. \quad (7.2.1)$$

(We use the letter W because such quantities are frequently termed *Wasserstein distances*, though we shall stick with our notation.) To understand these as a transportation cost, we think of the coupling π as a “plan” for moving probability mass from P to Q via some sampling scheme $\pi(\cdot | X)$ or $\pi(\cdot | Y)$.

In analogy with the total variation distance, which uses $c(x, y) = \mathbf{1}\{x \neq y\}$, the discrete distance, we can consider other costs on $\mathcal{X} \times \mathcal{Y}$. Notably, for any pair of functions $f : \mathcal{X} \rightarrow \mathbb{R}$ and $g : \mathcal{Y} \rightarrow \mathbb{R}$, whenever $f(x) + g(y) \leq c(x, y)$ for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$, we have the trivial inequality

$$\mathbb{E}_{\pi}[c(X, Y)] \geq \mathbb{E}_{\pi}[f(X) + g(Y)] = \mathbb{E}_P[f(X)] + \mathbb{E}_Q[g(Y)].$$

In fact, deep duality results for these distances. The next theorem, a variant of results typically called the *Kantorovich duality theorems*, captures the main results.

Theorem 7.2.1. *Let $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+ \cup \{\infty\}$ be a lower semicontinuous cost function and \mathcal{X} and \mathcal{Y} be metric spaces. Then for any distributions P and Q*

$$W_c(P, Q) = \sup \{\mathbb{E}_P[f(X)] + \mathbb{E}_Q[g(Y)] \mid f(x) + g(y) \leq c(x, y) \text{ for all } x \in \mathcal{X}, y \in \mathcal{Y}\},$$

and there is a coupling $\pi \in \Pi(P, Q)$ achieving $\mathbb{E}_{\pi}[c(X, Y)] = W_c(P, Q)$.

We provide a heuristic sketch of a proof of Theorem 7.2.1 in Section 7.2.2 to suggest why we might expect duality to hold, with pointers to more rigorous references. It is also typically the case that the supremum is attained by some functions f, g as well, but this is beyond our scope.

Most frequently, we consider cost functions c that are distances, meaning that $c : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$ satisfies the triangle inequality. In this case, we can restrict the supremum to be over only 1-Lipschitzian functions, where for a function f we define the Lipschitzian norm

$$\|f\|_{\text{Lip}, c} := \sup_{x \neq y} \frac{|f(x) - f(y)|}{c(x, y)}.$$

In this case, we have the following corollary, whose proof we defer to Section 7.2.3.

Corollary 7.2.2. *In addition to the conditions of Theorem 7.2.1, assume that c is a distance on $\mathcal{X} \times \mathcal{X}$. Then*

$$W_c(P, Q) = \min_{\pi \in \Pi(P, Q)} \mathbb{E}_{\pi}[c(X, Y)] = \sup \left\{ \mathbb{E}_P[f(X)] - \mathbb{E}_Q[f(X)] \mid \|f\|_{\text{Lip}, c} \leq 1 \right\}$$

For example, for the total variation distance, we have $c(x, y) = \mathbf{1}\{x \neq y\}$, and f being Lipschitzian is equivalent to $|f(x) - f(y)| \leq 1$ for all x, y .

7.2.1 A tensorized transportation inequality

The power of transportation inequalities comes from their ability to tensorize—one can extend a one-dimensional transportation inequality to an n -dimensional one using properties of the KL-divergence and convexity. We give one such argument here. Our starting point is a product distribution where each P_i satisfies a (modified) transportation cost inequality for an increasing convex function $\phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$:

$$\phi(W_c(Q, P_i)) \leq D_{\text{kl}}(Q \| P_i) \quad \text{for all } Q \quad (7.2.2)$$

for $i = 1, \dots, n$. For example, Pinsker's inequality gives the bound (7.2.2) for any distribution P_i with cost $c(x, y) = \mathbf{1}\{x \neq y\}$ and $\phi(t) = 2t^2$. The following theorem shows how we can leverage marginal transport

Theorem 7.2.3. *Let the distributions P_i on sets \mathcal{X}_i , $i = 1, \dots, n$, satisfy the marginal transportation bound (7.2.2). Then the product $P = P_1 \times \dots \times P_n$ satisfies*

$$\inf_{\pi \in \Pi(P, Q)} \sum_{i=1}^n \phi(\mathbb{E}_{\pi}[c(X_i, Y_i)]) \leq D_{\text{kl}}(Q \| P)$$

for all distributions Q on $\mathcal{X}_1 \times \dots \times \mathcal{X}_n$.

Proof We provide an inductive proof. Assume that the claimed inequality holds for some value n (the base case, of course, being the assumed inequality (7.2.2)). By the chain rule (2.1.6) for the KL-divergence, where we use the notation there as well for $X_1^{n+1} \sim Q$ and $Y_1^{n+1} \sim P$, we have

$$D_{\text{kl}}(Q \| P^{n+1}) = D_{\text{kl}}(X_1^n \| Y_1^n \mid X_{n+1}) + D_{\text{kl}}(X_{n+1} \| Y_{n+1}).$$

Now, let π_{n+1} be the optimal coupling between Q_{n+1} and P_{n+1} , and for values $x_{n+1} \in \mathcal{X}$, let π_x be the optimal coupling between $Q(X_1^n \in \cdot \mid X_{n+1} = x)$ and P^n (these exist by Theorem 7.2.1). Then by the induction hypothesis,

$$\sum_{i=1}^n \phi(\mathbb{E}_{\pi_x}[c(X_i, Y_i) \mid X_{n+1} = x]) \leq D_{\text{kl}}(X_1^n \| Y_1^n \mid X_{n+1} = x)$$

and

$$\phi(\mathbb{E}_{\pi_{n+1}}[c(X_{n+1}, Y_{n+1})]) \leq D_{\text{kl}}(X_{n+1} \| Y_{n+1}).$$

Then Jensen's inequality implies that, integrating over $X_{n+1} \sim Q_{n+1}$, we have

$$\begin{aligned} \sum_{i=1}^n \phi(\mathbb{E}_{\pi}[c(X_i, Y_i)]) + \phi(\mathbb{E}_{\pi_{n+1}}[c(X_{n+1}, Y_{n+1})]) \\ \leq D_{\text{kl}}(X_1^n \| Y_1^n \mid X_{n+1}) + D_{\text{kl}}(X_{n+1} \| Y_{n+1}) = D_{\text{kl}}(Q \| P^{n+1}). \end{aligned}$$

The distribution $\pi = \mathbb{E}_{Q_{n+1}}[\pi_{X_{n+1}}]$ and π_{n+1} have a consistent joint distribution by construction, giving the theorem. \square

In Section 7.3, we develop several applications of Theorem 7.2.3. Here, we present a few corollaries of the result.

Corollary 7.2.4 (Marton's transportation inequality). *Let P_i be distributions on \mathcal{X} and $P = P_1 \times \cdots \times P_n$. Then for any distribution Q on \mathcal{X}^n ,*

$$\inf_{\pi \in \Pi(P, Q)} \sum_{i=1}^n \pi(X_i \neq Y_i)^2 \leq \frac{1}{2} D_{\text{kl}}(Q \| P).$$

Proof For any Q and $Y \sim P_i$, we have by Pinsker's inequality and Theorem 7.2.1 that

$$\inf_{\pi \in \Pi(P_i, Q)} \pi(X \neq Y) = \|Q - P_i\|_{\text{TV}} \leq \sqrt{\frac{1}{2} D_{\text{kl}}(Q \| P_i)}.$$

This is the assumption of Theorem 7.2.3 with $\phi(t) = 2t^2$. □

We can recover the bounded differences inequality (Proposition 4.2.5) as a corollary (which, of course, we have already proved by easier means). To see this, note that any function $f : \mathcal{X}^n \rightarrow \mathbb{R}$ satisfying bounded differences with constants $b = [b_i]_{i=1}^n$ satisfies

$$|f(x_1^n) - f(y_1^n)| \leq \sum_{i=1}^n b_i \mathbf{1}\{x_i \neq y_i\}.$$

We use Corollary 7.2.4 coupled with Theorem 7.1.2 to prove the inequality. Letting $Z(x_1^n) = f(x_1^n) - \mathbb{E}_P[f(X_1^n)]$, we have for any coupling $\pi \in \Pi(Q, P)$ that

$$\mathbb{E}_Q[Z(Y_1^n)] - \mathbb{E}_P[Z(X_1^n)] \leq \sum_{i=1}^n b_i \pi(Y_i \neq X_i) \leq \left(\sum_{i=1}^n b_i^2 \right)^{1/2} \left(\sum_{i=1}^n \pi(Y_i \neq X_i)^2 \right)^{1/2}$$

by Cauchy-Schwarz. Applying Corollary 7.2.4, we obtain

$$\mathbb{E}_Q[Z] - \mathbb{E}_P[Z] \leq \|b\|_2 \sqrt{D_{\text{kl}}(Q \| P) / 2},$$

and so Theorem 7.1.2 gives

$$\mathbb{E}_P[e^{\lambda(f(X_1^n) - \mathbb{E}[f(X_1^n)])}] \leq \exp\left(\frac{\lambda^2 \|b\|_2^2}{8}\right)$$

for all λ , giving the same bound as Proposition 4.2.5.

7.2.2 A heuristic proof of Theorem 7.2.1

For a formal statement of this result, and much more detail, we refer to Villani [188]; Theorem 4.1 of his book guarantees the existence of an optimal coupling so long as c is lower semicontinuous, while Theorem 5.10 provides the duality guarantees.

Our proof is non-rigorous and mostly for intuition; to make the writing simpler, we shall assume P and Q have densities p and q for Lebesgue measure, and we will dispense with all measure-theoretic details and rigor. Letting π also be the density of (X, Y) , we write the initial transport problem as

$$\begin{aligned} & \text{minimize} && \mathbb{E}_\pi[c(X, Y)] \\ & \text{subject to} && \int_{\mathcal{Y}} \pi(x, y) dy = p(x) \quad \text{and} \quad \int_{\mathcal{X}} \pi(x, y) dx = q(y) \end{aligned}$$

over $\pi \geq 0$. Introduce Lagrange multiplier functions $f : \mathcal{X} \rightarrow \mathbb{R}$ and $g : \mathcal{Y} \rightarrow \mathbb{R}$ to obtain the Lagrangian

$$\begin{aligned} L(\pi, f, g) &= \int c(x, y) \pi(x, y) dx dy + \int_{\mathcal{X}} f(x) \left(p(x) - \int_{\mathcal{Y}} \pi(x, y) dy \right) dx + \int_{\mathcal{Y}} g(y) \left(q(y) - \int_{\mathcal{X}} \pi(x, y) dx \right) dy \\ &= \int (c(x, y) - f(x) - g(y)) \pi(x, y) dx dy + \mathbb{E}_P[f(X)] + \mathbb{E}_Q[g(Y)]. \end{aligned}$$

To obtain a dual problem, we minimize out the nonnegative function π . If there exists a pair (x, y) with $f(x) + g(y) > c(x, y)$ (again, eliding rigor), then we could send $\pi(x, y) \uparrow \infty$, making the first integral $-\infty$. Thus we obtain

$$\inf_{\pi \geq 0} L(\pi, f, g) = \begin{cases} \mathbb{E}_P[f(X)] + \mathbb{E}_Q[g(Y)] & \text{if } f(x) + g(y) \leq c(x, y) \text{ all } x \in \mathcal{X}, y \in \mathcal{Y} \\ -\infty & \text{otherwise.} \end{cases}$$

Assuming strong duality obtains, this gives the associated dual problem and theorem.

7.2.3 Proof of Corollary 7.2.2

Recognize that given any function f , we can only improve the objective by taking g as large as possible. Accordingly, we define the c -conjugate of f by

$$f^c(y) := \inf_x \{c(x, y) - f(x)\},$$

which at each y is the largest possible value v satisfying $f(x) + v \leq c(x, y)$ for all x . Setting $g = f^c$ can only increase the objective in the supremum. We therefore have the equivalent problem

$$\sup \{ \mathbb{E}_P[f(X)] + \mathbb{E}_Q[f^c(Y)] \},$$

where we note that $f(x) + f^c(y) \leq c(x, y)$ for all x, y . The conjugate f^c is Lipschitz with respect to the cost (distance) c : we have

$$f^c(x) - f^c(y) = \inf_{y'} \sup_{x'} \{c(x, y') - f(y') - c(x', y) + f(x')\} \leq \sup_{x'} \{c(x, x') - c(x', y)\} \leq c(x, y)$$

by the triangle inequality, and the lower bound is similar, showing that f^c is Lipschitz.

Of course, we should set f as large as possible given f^c while satisfying $f(x) + f^c(y) \leq c(x, y)$, meaning we choose

$$f^*(x) = \inf_y \{c(x, y) - f^c(y)\}.$$

But for this we note that because f^c is Lipschitz,

$$c(x, y) - f^c(y) \geq c(x, y) - f^c(x) - c(x, y) = -f^c(x),$$

and $\inf_y \{c(x, y) - f^c(y)\} = -f^c(x)$. That is, $f(x) = -f^c(x)$ and is therefore Lipschitz, giving the corollary.

7.3 Some applications of concentration and the variational inequality

In this section, we give two applications (though perhaps “applications” would be better nomenclature) for the transportation inequalities we have developed. The first is to concentration inequalities that depend more directly on the underlying geometry of a set \mathcal{X} and probability distribution P on \mathcal{X} .

7.3.1 Metric Gaussianity, transport inequalities, and expansion of sets

The connection between the transportation inequalities we have presented and concentration allows us to move beyond concentration of individual functions to more “geometric” concentration properties. The general formulation of such transportation inequalities, while beyond the scope of this book, allows us to develop concentration properties for many probability measures and types of functions. Here, we give one of the “basic” forms of this; in the next section we connect this to optimal testing, highlighting one of the ways that concentration and convergence interact with information theory. In Chapter 10.3 to come, we will present further applications of these results to fundamental limits in statistics and communication.

We begin by a general definition:

Definition 7.1. Let (\mathcal{X}, ρ) be a metric space with metric $\rho : \mathcal{X} \times \mathcal{X}$, and let $p \in [1, 2]$. Define the cost function $c_p(x, y) := \rho^p(x, y)$. A probability distribution P on \mathcal{X} satisfies an L^p transportation cost inequality with constant σ^2 , or a $T_p(\sigma^2)$ -inequality, if for all distributions Q on \mathcal{X}

$$W_{c_p}(Q, P)^{1/p} \leq \sqrt{2\sigma^2 D_{\text{kl}}(Q \| P)}.$$

We shall only consider $T_1(\sigma)$ -inequalities in this section, as they admit the most straightforward arguments. The bibliographic section provides pointers to other results.

Pinsker’s inequality shows if we use the Hamming metric on \mathcal{X} , then *all* distributions P satisfy a $T_1(\frac{1}{4})$ -inequality. Marton’s transportation inequality (Corollary 7.2.4) also shows that the product measure P^n satisfies a T_1 -inequality, but with a constant that depends on n : let d_{ham} be the Hamming distance on \mathcal{X}^n . Then for any coupling π on $X_1^n \sim Q$ and $Y_1^n \sim P^n$, we have

$$\mathbb{E}_\pi[d_{\text{ham}}(X_1^n, Y_1^n)] = \sum_{i=1}^n \pi(X_i \neq Y_i) \leq \sqrt{n} \left(\sum_{i=1}^n \pi(X_i \neq Y_i)^2 \right)^{1/2},$$

and so

$$W_{d_{\text{ham}}}(Q, P^n) \leq \sqrt{\frac{n}{2}} \sqrt{D_{\text{kl}}(Q \| P^n)},$$

a $T_1(\frac{n}{4})$ -inequality.

The astounding consequence of such transport inequalities is that probability measure P satisfying them then necessarily inherits some strong concentration properties. To discuss this we require a bit of additional notation. For a metric space (\mathcal{X}, ρ) and set $A \subset \mathcal{X}$, define the r -blowup or r -expansion of A by

$$A_r := \{x \in \mathcal{X} \mid \rho(x, A) \leq r\},$$

where we recall the notation $\rho(x, A) = \inf_{y \in A} \rho(x, y)$. When a small blowup of A drastically increases its probability, then P exhibits strong concentration properties. A more quantitative

version of this follows: we say P exhibits *Gaussian concentration on* (\mathcal{X}, ρ) with constant $\kappa > 0$ if for some $K < \infty$,

$$P(A) \geq \frac{1}{2} \text{ implies } P(A_r) \geq 1 - Ke^{-\kappa r^2} \text{ for } r \geq 0. \quad (7.3.1)$$

Because we often think of $r \uparrow \infty$, the only important constant is $\kappa > 0$, and so sometimes we provide the weaker guarantee $P(A_r) \geq 1 - Ke^{-\kappa[r-r_0]_+^2}$ for some $r_0 \geq 0$. This latter guarantee implies the Gaussian metric concentration (7.3.1) for large r with a worse constant K and (asymptotically) equivalent κ .

By a clever argument involving conditioning on belonging to a particular set A , it is possible to show the *blowing up* lemma: whenever a distribution P satisfies a transportation inequality, it enjoys Gaussian metric concentration.

Theorem 7.3.1 (The blowing up lemma). *Let P satisfy a $T_1(\sigma^2)$ -inequality. Then for any measurable set A with $P(A) > 0$ and $r \geq 0$,*

$$P(A) \cdot P(A_r^c) \leq \exp\left(-\frac{r^2}{4\sigma^2}\right)$$

and

$$P(A_r) \geq 1 - \exp\left(-\frac{1}{2\sigma^2} \left[r - \sqrt{2\sigma^2 \log \frac{1}{P(A)}}\right]_+^2\right).$$

Letting A be any set with $P(A) \geq \frac{1}{2}$, the first inequality in Theorem 7.3.1 implies that $P(A_r^c) \leq 2 \exp(-\frac{r^2}{4\sigma^2})$, that is,

$$P(A_r) \geq 1 - 2 \exp\left(-\frac{r^2}{4\sigma^2}\right).$$

The second inequality offers better constants in the exponent: taking $r_0 = \sqrt{2\sigma^2 \log 2}$ and $r = t + r_0$, we always have

$$P(A_{r_0+t}) \geq 1 - \exp\left(-\frac{t^2}{2\sigma^2}\right).$$

Proof For any set $A \subset \mathcal{X}$, define the conditional distribution

$$P_A(S) := P(A \cap S) / P(A).$$

Then without loss of generality, we may assume P_A and P have densities p_A and p with respect to some base measure μ , where $p_A(x) = p(x)/P(A)$. Then by inspection,

$$D_{\text{kl}}(P_A \| P) = \int p_A(x) \log \frac{p_A(x)}{p(x)} d\mu(x) = \log \frac{1}{P(A)}.$$

Then for any sets A, B , we have

$$W_\rho(P_A, P_B) \leq W_\rho(P_A, P) + W_\rho(P_B, P)$$

by the triangle inequality, and so

$$W_\rho(P_A, P_B) \leq \sqrt{2\sigma^2 \log \frac{1}{P(A)}} + \sqrt{2\sigma^2 \log \frac{1}{P(B)}}. \quad (7.3.2)$$

Now, for any joint $\pi \in \Pi(P_A, P_B)$ on $X \in A$ and $Y \in B$, we can write

$$\mathbb{E}_\pi[\rho(X, Y)] \geq \inf_{x \in A} \inf_{y \in B} \rho(x, y) = \rho(A, B). \quad (7.3.3)$$

Substituting this into the bound (7.3.2), we have for any (measurable) sets A, B with $P(A) > 0$ and $P(B) > 0$ that

$$\rho(A, B) \leq \sqrt{2\sigma^2 \log \frac{1}{P(A)}} + \sqrt{2\sigma^2 \log \frac{1}{P(B)}}.$$

From here, we can perform two manipulations. The first is to recognize that by the inequality $(\sqrt{a} + \sqrt{b})^2 = a + b + 2\sqrt{ab} \leq 2a + 2b = (\sqrt{2(a+b)})^2$, we have $\sqrt{a} + \sqrt{b} \leq \sqrt{2}\sqrt{a+b}$, and so

$$\rho(A, B) \leq \sqrt{4\sigma^2 \log \frac{1}{P(A)} \frac{1}{P(B)}}.$$

Now if we take B be the set of points at least r far away: recalling the r -blowup $A_r := \{y \in \mathcal{X} \mid d(y, A) \leq r\}$ of A , we set $B = A_r^c$ and obtain $\rho(A, B) = \rho(A, A_r^c) \geq r$ for all $r \geq 0$, and $P(B) = P(A_r^c) = 1 - P(A_r)$. In particular,

$$r \leq 2\sqrt{\sigma^2 \log \frac{1}{P(A)P(A_r^c)}}.$$

Rearranging gives $P(A)P(A_r^c) \leq \exp(-\frac{r^2}{4\sigma^2})$. Alternatively, we can use inequality (7.3.2) directly to obtain

$$r \leq \sqrt{2\sigma^2 \log \frac{1}{P(A)}} + \sqrt{2\sigma^2 \log \frac{1}{1 - P(A_r)}}.$$

Rearranging and solving for $P(A_r)$ gives the theorem. \square

Theorem 7.3.1 is particularly evocative for product measures on discrete spaces. In this case, we use the Hamming metric d_{ham} , which satisfies a $T_1(\frac{n}{4})$ -transport cost inequality (Definition 7.1).

Corollary 7.3.2. *Let $A \subset \mathcal{X}^n$ be any measurable set and let $X_i \stackrel{\text{iid}}{\sim} P$. Then for any $r \geq 0$,*

$$\mathbb{P}(d_{\text{ham}}(X_1^n, A) \leq r) \geq 1 - \exp\left(-\frac{2}{n} \left[r - \sqrt{\frac{n}{2} \log \frac{1}{P^n(A)}}\right]_+^2\right)$$

As a consequence, for any sequence of sets that have at least polynomial probability—the sets $A_n \subset \mathcal{X}^n$ satisfy $P^n(A_n) \geq n^{-p}$ for some $p < \infty$ —neighborhoods expanded by only \sqrt{n} in the Hamming metric necessarily cover nearly the entire probability mass. Indeed, in this case, $\log \frac{1}{P^n(A)} \leq p \log n$, and so setting

$$r_n = c\sqrt{\frac{n}{2}} + \sqrt{\frac{pn}{4} \log n},$$

we have

$$\mathbb{P}(d_{\text{ham}}(X_1^n, A_n) \leq r_n) \geq 1 - \exp(-c^2) \text{ for all } n.$$

7.3.2 A weak and strong converse for hypothesis testing

Let us return to the simplest statistical question that underlies much of our development: a simple hypothesis test. In this case, we wish to test a null H_0 against an alternative H_1 ,

$$H_0 : X_1^n \stackrel{\text{iid}}{\sim} P_0 \text{ versus } H_1 : X_1^n \stackrel{\text{iid}}{\sim} P_1,$$

where P_0 and P_1 are different distributions. Chapter 2.3.1 discusses this simple problem, connecting bounds on the variation distance to lower bounds on the summed probability of error. Here, we give a different interpretation, fixing both P and Q and showing that if the type-I probability of error is small, the rate at which the type-II probability of error can only decrease at a certain rate.

To give intuition, let us first develop the optimal test, which gives an achievability result and shows the correct behavior of the asymptotic error; we provide the converse results after this. We shall assume for simplicity that the variance of the log-likelihood ratio $v^2 := \text{Var}_P(\log \frac{p_0(X)}{p_1(X)}) < \infty$. Define the log-likelihood ratio of a sequence x_1^n by

$$L_n(x_1^n) := \sum_{i=1}^n \log \frac{p_0(x_i)}{p_1(x_i)},$$

so that if $\Phi(t) = \mathbb{P}(Z \leq t)$ is the standard normal CDF, the central limit theorem implies

$$P_0^n(L_n(X_1^n) \geq nD_{\text{kl}}(P_0\|P_1) + v\sqrt{n}\Phi^{-1}(1 - \epsilon)) \rightarrow \epsilon.$$

The Neyman-Pearson lemma gives that the optimal test of P_0 against P_1 is to compare $L_n(X_1^n)$ against a threshold, and so we may take a sequence $a_n(\epsilon)$ for which $a_n(\epsilon)/\sqrt{n} \rightarrow 0$ and define the test

$$\Psi_n(x_1^n) = \begin{cases} \text{accept } H_0 & \text{if } L_n(x_1^n) - nD_{\text{kl}}(P_0\|P_1) \geq v\sqrt{n} + a_n(\epsilon) \\ \text{accept } H_1 & \text{otherwise.} \end{cases}$$

Letting $\Psi = 0$ indicate accepting H_0 and $\Psi = 1$ indicate accepting H_1 for simplicity, this test satisfies $P_0^n(\Psi_n = 0) \geq 1 - \epsilon$ for all n ; taking $a_n(\epsilon)$ as large as possible while satisfying this guarantee yields an optimal test (via the Neyman-Pearson lemma) with $a_n(\epsilon) = o(\sqrt{n})$.

Now we consider the type-II error, that is, $P_1^n(\Psi_n = 0)$. For this, we observe that

$$\begin{aligned} P_1^n(\Psi_n = 0) &= P_1^n(L_n(X_1^n) - nD_{\text{kl}}(P_0\|P_1) \geq v\sqrt{n}\Phi^{-1}(1 - \epsilon) + a_n(\epsilon)) \\ &\stackrel{(i)}{\leq} \mathbb{E}_{P_1^n} \left[e^{L_n(X_1^n)} \right] \exp(-nD_{\text{kl}}(P_0\|P_1) - v\sqrt{n}\Phi^{-1}(1 - \epsilon) - a_n(\epsilon)) \\ &\stackrel{(ii)}{=} \exp(-nD_{\text{kl}}(P_0\|P_1) - v\sqrt{n}\Phi^{-1}(1 - \epsilon) - a_n(\epsilon)), \end{aligned}$$

where inequality (i) is a Chernoff bound and equality (ii) holds because $\mathbb{E}_{P_1}[\exp(\log \frac{p_0(X)}{p_1(X)})] = 1$. Summarizing, we have the following result.

Proposition 7.3.3. *Let $\epsilon \in (0, 1)$ and let Ψ_n be the level- ϵ likelihood ratio test of H_0 against H_1 , that is, $P^n(\Psi_n = H_1) \leq \epsilon$, and assume that $v^2 := \text{Var}_P(\log \frac{p(X)}{q(X)}) < \infty$. Then*

$$\frac{1}{n} \log \frac{1}{P_1^n(\Psi_n = 0)} \leq -D_{\text{kl}}(P_0\|P_1) - \frac{v}{\sqrt{n}}\Phi^{-1}(1 - \epsilon) + o(1/\sqrt{n}).$$

We show now that we can derive converse results providing the similar rates of convergence using the blowing-up lemma (Theorem 7.3.1) without relying so carefully on the particulars of the simple hypothesis test, especially the Neyman-Pearson lemma. We first present a so-called “weak converse”, then extend the idea to a “strong converse.” One should think of this distinction as follows: a weak converse states that if the probability of (type-I) error tends to zero, then the asymptotic type-II error can only decrease to zero at a certain rate. This does not, however, eliminate the possibility of lower type-II errors by allowing a fixed $\epsilon > 0$ probability of error in testing the null H_0 . The strong converse, on the other hand, eliminates this possibility: for any fixed $\epsilon > 0$, *no* improvement in (asymptotic) type-II error is possible. We note that the style of argument we employ here extends beyond this simple setting: Chapter 10.3 employs it to give similar fundamental limits in some communication and estimation problems.

Both the weak and strong converse we prove follow from the same simple idea: the data processing inequality for the KL-divergence.

Proposition 7.3.4 (Weak converse in hypothesis testing). *Let P_0 and P_1 be arbitrary distributions, and let the test Ψ have level ϵ , that is, $P_0(\Psi = 1) \leq \epsilon$. Then*

$$-\log \frac{1}{P_1(\Psi = 0)} \geq -\frac{D_{\text{kl}}(P_0 \| P_1)}{1 - \epsilon} - \frac{\log 2}{1 - \epsilon}.$$

Specializing the result to the product distribution case, replacing P_0 and P_1 with P_0^n and P_1^n , respectively, we obtain that

$$-\log \frac{1}{P_1^n(\Psi_n = 0)} \geq -\frac{n D_{\text{kl}}(P_0 \| P_1)}{1 - \epsilon} - \frac{\log 2}{1 - \epsilon},$$

and so for any $\epsilon_n \rightarrow 0$, we have the limit

$$\liminf_{n \rightarrow \infty} -\frac{1}{n} \log \frac{1}{P_1^n(\Psi_n = 0)} \geq -D_{\text{kl}}(P_0 \| P_1),$$

which matches Proposition 7.3.3.

Proof Let $A = \{x \in \mathcal{X} \mid \Psi(A) = 0\}$, and let $p_0 = P_0(A)$ and $p_1 = P_1(A)$. Recall our notation for the binary relative entropy $D_{\text{kl}}(p \| q) = p \log \frac{p}{q} + (1 - p) \log \frac{1-p}{1-q}$ and the binary entropy functional $h_2(p) = -p \log p - (1 - p) \log(1 - p) \leq \log 2$. Then by the data-processing inequality

$$\begin{aligned} D_{\text{kl}}(P_0 \| P_1) &\geq D_{\text{kl}}(P_0(A) \| P_1(A)) = p_0 \log \frac{p_0}{p_1} + (1 - p_0) \log \frac{1 - p_0}{1 - p_1} \\ &= p_0 \log \frac{1}{p_1} + (1 - p_0) \log \frac{1}{1 - p_1} - h_2(p_0) \\ &\geq p_0 \log \frac{1}{p_1} - h_2(p_0). \end{aligned}$$

Rearranging gives $\log \frac{1}{p_1} \geq -\frac{1}{p_0} D_{\text{kl}}(P_0 \| P_1) + \frac{\log 2}{p_0}$, and substituting $p_0 \geq 1 - \epsilon$ gives the result. \square

Let us now provide the stronger result. An unfortunate limitation of the method we employ here is that we must assume \mathcal{X} is finite to enable the strongest applications of the blowing up lemma (Theorem 7.3.1). In this case, we focus on the product distributions P_0^n and P_1^n , and let \mathcal{X} be the support of P_0 . We assume without loss of generality that $\inf_{x \in \mathcal{X}} P_1(\{x\}) > 0$, as otherwise,

$D_{\text{kl}}(P_0\|P_1) = \infty$ and any lower bound is trivial. The key insight, which Ahlswede et al. [4] develop, is that we should apply the data processing inequality to the enlargement $A_r = \{x_1^n \mid d_{\text{ham}}(x_1^n, A) \leq r\}$ of the acceptance set $A = \{x \in \mathcal{X}^n \mid \Psi_n(x) = 0\}$ rather than A itself; because of the blowing-up lemma, this set has nearly full measure, which in turn means that $P_1(\Psi_n = 0)$ cannot actually be too small.

Proposition 7.3.5 (Strong converse in hypothesis testing). *Let P_0 and P_1 be distributions on \mathcal{X} , and let the test Ψ_n have level ϵ for testing P_0^n against P_1^n . Then*

$$\frac{1}{n} \log \frac{1}{P_1^n(\Psi_n = 0)} \geq -D_{\text{kl}}(P_0\|P_1) - O(1) \frac{\log n \sqrt{\log \frac{n}{1-\epsilon}}}{\sqrt{n}}$$

for all large enough n .

Proof Following the paragraph preceding the proposition, let $A = \{x_1^n \in \mathcal{X}^n \mid \Psi_n(x_1^n) = 0\}$ be the acceptance set of Ψ_n , and let $A_r = \{x \mid d_{\text{ham}}(x, A) \leq r\}$ its r -blowup. For simplicity in notation, define $p_r = P_0^n(A_r)$ and $q_r = P_1^n(A_r)$. Then by the data-processing inequality as in the proof of Proposition 7.3.4, we have

$$nD_{\text{kl}}(P_0\|P_1) = D_{\text{kl}}(P_0^n\|P_1^n) \geq D_{\text{kl}}(P_0^n(A_r)\|P_1^n(A_r)) = D_{\text{kl}}(p_r\|q_r) \geq p_r \log \frac{1}{q_r} - h_2(p_r). \quad (7.3.4)$$

Of course, we wish to provide bounds on $\log \frac{1}{q}$ for $q = P_1^n(A)$, not q_r . For this, we use the following combinatorial lemma:

Lemma 7.3.6. *Define $q_\star = \min_{x \in \mathcal{X}} P_1(\{x\})$. For any r ,*

$$P_1^n(A) \leq P_1^n(A_r) \leq \binom{n}{r} \left(\frac{\text{card}(\mathcal{X})}{q_\star} \right)^r P_1^n(A).$$

Proof Let $x, x' \in \mathcal{X}^n$. Then

$$P_1^n(\{x\}) = \prod_{i: x_i = x'_i} P_1(\{x_i\}) \prod_{i: x_i \neq x'_i} P_1(\{x_i\}) = P_1(\{x'\}) \prod_{i: x_i \neq x'_i} \frac{P_1(\{x_i\})}{P_1(\{x'_i\})} \geq P_1(\{x'\}) q_\star^{d_{\text{ham}}(x, x')}.$$

So for any single $x \in \mathcal{X}^n$ and its r -neighborhood $\{x\}_r = \{x' \mid d_{\text{ham}}(x, x') \leq r\}$, we have

$$P_1^n(\{x\}_r) \leq \binom{n}{r} \left(\frac{\text{card}(\mathcal{X})}{q_\star} \right)^r P_1^n(\{x\}).$$

The union bound implies the lemma. \square

By the lemma, we therefore obtain for the P_1 -dependent constant $K = \frac{\text{card}(\mathcal{X})}{q_\star}$ and $q := P_1^n(A) = P_1^n(\Psi_n = 0)$, we have

$$q_r \leq \binom{n}{r} K^r q \leq \left(\frac{Kne}{r} \right)^r q$$

by the standard binomial bound $\binom{n}{r} \leq \left(\frac{ne}{r} \right)^r$. Substituting into the data processing bound (7.3.4), we obtain

$$nD_{\text{kl}}(P_0\|P_1) \geq p_r \log \frac{1}{q} - r \log \frac{Kne}{r} - \log 2.$$

Finally, we use the blowing-up lemma (Theorem 7.3.1): Fix $t \geq 0$ to be chosen and take

$$r = r(t) := t\sqrt{\frac{n}{2}} + \sqrt{\frac{n}{2} \log \frac{1}{1-\epsilon}}.$$

Then $p_r \geq 1 - e^{-t^2}$, and so

$$nD_{\text{kl}}(P_0 \| P_1) \geq (1 - e^{-t^2}) \log \frac{1}{q} - r \log \frac{Kne}{r} - \log 2.$$

If we take $t = \sqrt{2 \log n}$ then $r = O(1)\sqrt{n \log \frac{n}{1-\epsilon}}$, and so

$$-(1 - n^{-2}) \log \frac{1}{q} \geq -nD_{\text{kl}}(P_0 \| P_1) - O(1)\sqrt{n \log^2 n \cdot \log \frac{n}{1-\epsilon}} + n \log n \cdot \log K.$$

Dividing by $n(1 - 1/n^2)$ gives the result. \square

7.4 Discussion and bibliographic remarks

JCD Comment: Might want to investigate a bit of the treatment in both Raginsky and Sason [158] and Liu et al. [139].

JCD Comment: Maybe add some refs to the paper [199], and give a citation for Corollary 7.1.4 to Catoni and Giulini [49] or whomever is appropriate.

7.5 Exercises

Exercise 7.1: Complete the proof of Corollary 7.1.1 in the case that $\mathbb{E}_Q[e^{g(X)}] = +\infty$.

Exercise 7.2 (A discrete isoperimetric inequality): Let $A \subset \mathbb{Z}^d$ be a finite subset of the d -dimensional integers. Let the projection mapping $\pi_j : \mathbb{Z}^d \rightarrow \mathbb{Z}^{d-1}$ be defined by

$$\pi_j(z_1, \dots, z_d) = (z_1, \dots, z_{j-1}, z_{j+1}, \dots, z_d)$$

so that we “project out” the j th coordinate, and define the projected sets

$$\begin{aligned} A_j &= \pi_j(A) = \{\pi_j(z) : z \in A\} \\ &= \left\{ z \in \mathbb{Z}^{d-1} : \text{there exists } z_\star \in \mathbb{Z} \text{ such that } (z_1, z_2, \dots, z_{j-1}, z_\star, z_j, \dots, z_{d-1}) \in A \right\}. \end{aligned}$$

Prove the Loomis-Whitney inequality, that is, that

$$\text{card}(A) \leq \left(\prod_{j=1}^d \text{card}(A_j) \right)^{\frac{1}{d-1}}.$$

Exercise 7.3: Let X be a non-constant and mean-zero random variable with moment generating function φ_X defined on a neighborhood of 0, and let

$$b_1 = \sup \{ \lambda \geq 0 \mid \varphi_X(\lambda) < \infty \} \quad \text{and} \quad b_0 = \inf \{ \lambda \leq 0 \mid \varphi_X(\lambda) < \infty \}.$$

You may assume that φ_X is \mathcal{C}^∞ on (b_0, b_1) (Proposition 3.2.2 guarantees this).

- (a) Show that φ_X is strictly increasing on $(0, b_1)$ and strictly decreasing on $(b_0, 0)$.
- (b) Show that φ'_X is strictly increasing on (b_0, b_1) .
- (c) Show that $\varphi_X(\lambda) \rightarrow \varphi_X(b_1)$ as $\lambda \uparrow b_1$ and $\varphi_X(\lambda) \rightarrow \varphi_X(b_0)$ as $\lambda \downarrow b_0$. (These limits may be finite or infinite.)
- (d) Show that $\varphi'_X(\lambda) \rightarrow \mathbb{E}[Xe^{b_1 X}] > 0$ as $\lambda \uparrow b_1$. (This limit may be finite or infinite.)
- (e) Show that the domain of φ_X , even when non-trivial, may be open or closed. In particular, give an example of a random variable for which $\text{dom } \varphi_X = [-1, 1]$, and give an example of a random variable for which $\text{dom } \varphi_X = (-1, 1)$.

Exercise 7.4 (“Dimension free” covariance estimation): Define the *effective rank* of a matrix $\Sigma \succeq 0$ by

$$\text{r}_{\text{eff}}(\Sigma) := \frac{\text{tr}(\Sigma)}{\|\Sigma\|_{\text{op}}}.$$

Let $M_i \in \mathbb{R}^{d \times d}$ be independent positive definite matrices for which $\mathbb{E}[M_i] = \Sigma$, and assume

$$\|u^\top M_i u\|_{\psi_1} \leq \kappa^2 u^\top \Sigma u$$

for any $u \in \mathbb{R}^d$. This problem uses Corollary 7.1.4 to extend Proposition 7.1.5 to show that such random matrices concentrate similarly to isotropic random matrices, but the effective rank replaces the dimension d .

- (a) Fix $\beta > 0$ and $r > 0$, and define the prior distribution π_0 on $\Theta = \mathbb{R}^d \times \mathbb{R}^d$ by $\mathbf{N}(0, \beta^{-1}\Sigma) \times \mathbf{N}(0, \beta^{-1}\Sigma)$, that is, the product of two mean-zero normals with covariance $\beta^{-1}\Sigma$. Let P_u be the normal distribution $\mathbf{N}(u, \beta^{-1}\Sigma)$ truncated restricted to the ball $u + r\mathbb{B}_2^d = \{x \mid \|x - u\|_2 \leq r\}$, and define the posterior $\pi_{u,v} = P_u \times P_v$. Show that for any $u, v \in \Sigma^{1/2}\mathbb{S}^{d-1} = \{u \mid u^\top \Sigma^{-1}u = 1\}$, we have

$$D_{\text{kl}}(\pi_{u,v} \parallel \pi_0) = 2 \log \frac{1}{C(r)} + \beta,$$

where $C(r) := \mathbb{P}(\|Z\|_2 \leq r)$ for $Z \sim \mathbf{N}(0, \beta^{-1}\Sigma)$ is the normalization for P_u .

- (b) Show that if $r = \sqrt{2\beta^{-1} \text{tr}(\Sigma)}$ then $C(r) \geq \frac{1}{2}$, and conclude that $D_{\text{kl}}(\pi_{u,v} \parallel \pi_0) \leq 2 \log 2 + \beta$.
- (c) Define

$$f(\theta_1, \theta_2, M) := \theta_1^\top \Sigma^{-1/2} (M - \Sigma) \Sigma^{-1/2} \theta_2.$$

Show that for fixed θ_1, θ_2 , f is sub-exponential, specifically,

$$\|f(\theta_1, \theta_2, M)\|_{\psi_1} \leq \kappa^2 (\|\theta_1\|_2^2 + \|\theta_2\|_2^2).$$

- (d) Show that for $(\theta_1, \theta_2) \sim \pi_{u,v}$ with the choice $r = \sqrt{2 \text{tr}(\Sigma)/\beta}$, we have

$$\max \{\|\theta_1\|_2, \|\theta_2\|_2\} \leq \sqrt{\|\Sigma\|_{\text{op}}} + \sqrt{2\beta^{-1} \text{tr}(\Sigma)}$$

whenever $u, v \in \Sigma^{1/2}\mathbb{S}^{d-1}$.

- (e) Use Corollary 7.1.4 to show that for any fixed $|\lambda| \leq \frac{1}{16\kappa^2\|\Sigma\|_{\text{op}}}$, with probability at least $1 - t$,

$$\lambda u^\top \Sigma^{-1/2} (P_n M - \Sigma) \Sigma^{-1/2} v \leq 64\lambda^2 \kappa^4 \|\Sigma\|_{\text{op}} + \frac{2 \log 2 + 2\text{r}_{\text{eff}}(\Sigma) + t}{n}$$

for all vectors $u, v \in \Sigma^{1/2}\mathbb{S}^{d-1}$ by appropriate choice of β .

- (f) Show that so long as $n \geq 4\text{r}_{\text{eff}}(\Sigma) + t$, then

$$\|P_n M - \Sigma\|_{\text{op}} \leq 20\kappa^2 \|\Sigma\|_{\text{op}} \sqrt{\frac{4\text{r}_{\text{eff}}(\Sigma) + t}{n}}. \quad (7.5.1)$$

(The constants are unimportant.)

Exercise 7.5: In this question, we compare the convergence guarantee in Proposition 7.1.5 to that in inequality (7.5.1). Let $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \Sigma)$ for some covariance $\Sigma \succeq 0$.

- (a) Give the smallest value σ^2 for which

$$\|\langle X, u \rangle\|_{\psi_2} \leq \sigma \|u\|_2$$

for all $u \in \mathbb{R}^d$.

- (b) Give the smallest value κ^2 for which

$$\|\langle X, u \rangle\|_{\psi_2} \leq \kappa \sqrt{u^\top \Sigma u}$$

for all $u \in \mathbb{R}^d$.

- (c) Given the previous parts, compare the guarantees Proposition 7.1.5 and inequality (7.5.1) provide on the deviations $P_n X X^\top - \Sigma$. How far apart can they be?

Exercise 7.6 (An error lower bound): Let P_0 and P_1 be distributions on \mathcal{X} , where \mathcal{X} is finite. Let $\Psi_n : \mathcal{X}^n \rightarrow \{0, 1\}$ be a sequence of tests with small enough type-II error that

$$\gamma_n := \frac{1}{n} \log \frac{1}{P_1^n(\Psi_n = 0)} - D_{\text{kl}}(P_0 \| P_1) \geq 0$$

for all large enough n . Show that there is a numerical constant $c > 0$ such that for all large enough n , the type-I error $\epsilon := P_0^n(\Psi_n = 1)$ satisfies

$$\epsilon \geq 1 - n \exp\left(-c \frac{n\gamma_n^2}{\log^2 n}\right).$$

JCD Comment: Exercise on self-concordant losses and measuring error in the “natural” metric based on the Hessian.

Chapter 8

Privacy and disclosure limitation

In this chapter, we continue to build on our ideas on stability in different scenarios, ranging from model fitting and concentration to interactive data analyses. Here, we show how stability ideas allow us to provide a new type of protection: the privacy of participants in studies. Until the mid-2000s, the major challenge in this direction had been a satisfactory definition of privacy, because collection of side information often results in unforeseen compromises of private information. The introduction of *differential privacy*—a type of stability in likelihood ratios for data releases from differing samples—alleviated these challenges, providing a firm foundation on which to build private estimators and other methodology. (Though it is possible to trace some of the definitions and major insights in privacy back at least to survey sampling literature in the 1960s.) Consequently, in this chapter we focus on privacy notions based on differential privacy and its cousins, developing the information-theoretic stability ideas helpful to understand the protections it is possible to provide.

8.1 Disclosure limitation, privacy, and definitions

We begin this chapter with a few cautionary tales and examples, which motivate the coming definitions of privacy that we consider. A natural belief might be that, given only certain summary statistics of a large dataset, individuals in the data are protected. Yet this appears, by and large, to be false. As an example, in 2008 Nils Homer and colleagues [113] showed that even releasing aggregated genetic frequency statistics (e.g., frequency of single nucleotide polymorphisms (SNP) in microarrays) can allow resolution of individuals within a database. Consequently, the US National Institutes of Health (NIH), the Wellcome Trust, and the Broad Institute removed genetic summaries from public access (along with imposing stricter requirements for private access) [171, 56].

Another hypothetical example may elucidate some of the additional challenges. Suppose that I release a dataset that consists of the frequent times that posts are made worldwide that denigrate government policies, but I am sure to remove all information such as IP addresses, usernames, or other metadata excepting the time of the post. This might seem *a priori* reasonably safe, but now suppose that an authoritarian government knows precisely when its citizens are online. Then by linking the two datasets, the government may be able to track those who post derogatory statements about their leaders.

Perhaps the strongest definition of privacy of databases and datasets is due to Dalenius [62], who suggests that “nothing about an individual should be learnable from the database that cannot be learned without access to the database.” But quickly, one can see that it is essentially impossible to reconcile this idea with scientific advancement. Consider, for example, a situation where we

perform a study on smoking, and discover that smoking causes cancer. We publish the result, but now we have “compromised” the privacy of everyone who smokes who did not participate in the study: we know they are more likely to get cancer.

In each of these cases, the biggest challenge is one of side information: how can we be sure that, when releasing a particular statistic, dataset, or other quantity that no adversary will be able to infer sensitive data about participants in our study? We articulate three desiderata that—we believe—suffice for satisfactory definitions of privacy. In discussion of private releases of data, we require a bit of vocabulary. We term a (randomized) algorithm releasing data either a *privacy mechanism*, consistent with much of the literature in privacy, or a *channel*, mapping from the input sample to some output space, in keeping with our statistical and information-theoretic focus. In no particular order, we wish our privacy mechanism, which takes as input a sample $X_1^n \in \mathcal{X}^n$ and releases some Z to satisfy the following.

- i. Given the output Z , even an adversary knowing everyone in the study (excepting one person) should not be able to test whether you belong to the study.
- ii. If you participate in multiple “private” studies, there should be some graceful degradation in the privacy protections, rather than a catastrophic failure. As part of this, any definition should guarantee that further processing of the output Z of a private mechanism $X_1^n \rightarrow Z$, in the form of the Markov chain $X_1^n \rightarrow Z \rightarrow Y$, should not allow further compromise of privacy (that is, a data-processing inequality). Additional participation in “private” studies should continue to provide little additional information.
- iii. The mechanism $X_1^n \rightarrow Z$ should be resilient to side information: even if someone knows something about you, he should learn little about you if you belong to X_1^n , and this should remain true even if the adversary later gleans more information about you.

The third desideratum is perhaps most elegantly phrased via a Bayesian perspective, where an adversary has some prior beliefs π on the membership of a dataset (these prior beliefs can then capture any side information the adversary has). The strongest adversary has a prior supported on two samples $\{x_1, \dots, x_n\}$ and $\{x'_1, \dots, x'_n\}$ differing in only a single element; a private mechanism would then guarantee the adversary’s posterior beliefs (after the release $X_1^n \rightarrow Z$) should not change significantly.

Before continuing addressing these challenges, we take a brief detour to establish notation for the remainder of the chapter. It will be convenient to consider randomized procedures acting on samples themselves; a sample x_1^n is clearly isomorphic to the empirical distribution $P_n = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i}$, and for two empirical distributions P_n and P'_n supported on $\{x_1, \dots, x_n\}$ and $\{x'_1, \dots, x'_n\}$, we evidently have

$$n \|P_n - P'_n\|_{\text{TV}} = d_{\text{ham}}(\{x_1, \dots, x_n\}, \{x'_1, \dots, x'_n\}),$$

and so we will identify samples with their empirical distributions. With this notational convenience in place, we then identify

$$\mathcal{P}_n = \left\{ P_n = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i} \mid x_i \in \mathcal{X} \right\}$$

as the set of all empirical distributions on n points in \mathcal{X} and we also abuse notation in an obvious way to define $d_{\text{ham}}(P_n, P'_n) := n \|P_n - P'_n\|_{\text{TV}}$ as the number of differing observations in the samples P_n and P'_n represent. A mechanism M is then a (typically) randomized mapping $M : \mathcal{P}_n \rightarrow \mathcal{Z}$,

which we can identify with its induced Markov channel Q from $\mathcal{X}^n \rightarrow \mathcal{Z}$; we use the equivalent views as is convenient.

The challenges of side information motivate [Dwork et al.](#)'s definition of *differential privacy* [80]. The key in differential privacy is that the noisy channel releasing statistics provides guarantees of bounded likelihood ratios between neighboring samples, that is, samples differing in only a single entry.

Definition 8.1 (Differential privacy). *Let $M : \mathcal{P}_n \rightarrow \mathcal{Z}$ be a randomized mapping. Then M is ε -differentially private if for all (measurable) sets $S \subset \mathcal{Z}$ and all $P_n, P'_n \in \mathcal{P}_n$ with $d_{\text{ham}}(P_n, P'_n) \leq 1$,*

$$\frac{\mathbb{P}(M(P_n) \in S)}{\mathbb{P}(M(P'_n) \in S)} \leq e^\varepsilon. \quad (8.1.1)$$

The intuition and original motivation for this definition are that an individual has little incentive to participate (or not participate) in a study, as the individual's data has limited effect on the outcome.

The model (8.1.1) of differential privacy presumes that there is a trusted curator, such as a hospital, researcher, or corporation, who can collect all the data into one centralized location, and it is consequently known as the *centralized* model. A stronger model of privacy is the *local model*, in which data providers trust no one, not even the data collector, and privatize their individual data before the collector even sees it.

Definition 8.2 (Local differential privacy). *A channel Q from \mathcal{X} to \mathcal{Z} is ε -locally differentially private if for all measurable $S \subset \mathcal{Z}$ and all $x, x' \in \mathcal{X}$,*

$$\frac{Q(Z \in S \mid x)}{Q(Z \in S \mid x')} \leq e^\varepsilon. \quad (8.1.2)$$

It is clear that Definition 8.2 and the condition (8.1.2) are stronger than Definition 8.1: when samples $\{x_1, \dots, x_n\}$ and $\{x'_1, \dots, x'_n\}$ differ in at most one observation, then the local model (8.1.2) guarantees that the densities

$$\frac{dQ(Z_1^n \mid \{x_i\})}{dQ(Z_1^n \mid \{x'_i\})} = \prod_{i=1}^n \frac{dQ(Z_i \mid x_i)}{dQ(Z_i \mid x'_i)} \leq e^\varepsilon,$$

where the inequality follows because only a single ratio may contain $x_i \neq x'_i$.

In the remainder of this introductory section, we provide a few of the basic mechanisms in use in differential privacy, then discuss its “semantics,” that is, its connections to the three desiderata we outline above. In the coming sections, we revisit a few more advanced topics, in particular, the composition of multiple private mechanisms and a few weakenings of differential privacy, as well as more sophisticated examples.

8.1.1 Basic mechanisms

The basic mechanisms in either the local or centralized models of differential privacy use some type of noise addition to ensure privacy. We begin with the simplest and oldest mechanism, randomized response, for local privacy, due to Warner [190] in 1965.

Example 8.1.1 (Randomized response): We wish to have a participant in a study answer a yes/no question about a sensitive topic (for example, drug use). That is, we would like to

estimate the proportion of the population with a characteristic (versus those without); call these groups 0 and 1. Rather than ask the participant to answer the question specifically, however, we give them a spinner with a face painted in two known areas, where the first corresponds to group 0 and has area $e^\varepsilon/(1 + e^\varepsilon)$ and the second to group 1 and has area $1/(1 + e^\varepsilon)$. Thus, when the participant spins the spinner, it lands in group 0 with probability $e^\varepsilon/(1 + e^\varepsilon)$. Then we simply ask the participant, upon spinning the spinner, to answer “Yes” if he or she belongs to the indicated group, “No” otherwise.

Let us demonstrate that this randomized response mechanism provides ε -local differential privacy. Indeed, we have

$$\frac{Q(\text{Yes} \mid x = 0)}{Q(\text{Yes} \mid x = 1)} = e^{-\varepsilon} \quad \text{and} \quad \frac{Q(\text{No} \mid x = 0)}{Q(\text{No} \mid x = 1)} = e^\varepsilon,$$

so that $Q(Z = z \mid x)/Q(Z = z \mid x') \in [e^{-\varepsilon}, e^\varepsilon]$ for all x, z . That is, the randomized response channel provides ε -local privacy. \diamond

The interesting question is, of course, whether we can still use this channel to estimate the proportion of the population with the sensitive characteristic. Indeed, we can. We can provide a somewhat more general analysis, however, which we now do so that we can give a complete example.

Example 8.1.2 (Randomized response, continued): Suppose that we have an attribute of interest, x , taking the values $x \in \{1, \dots, k\}$. Then we consider the channel (of Z drawn conditional on x)

$$Z = \begin{cases} x & \text{with probability } \frac{e^\varepsilon}{k-1+e^\varepsilon} \\ \text{Uniform}([k] \setminus \{x\}) & \text{with probability } \frac{k-1}{k-1+e^\varepsilon}. \end{cases}$$

This (generalized) randomized response mechanism is evidently ε -locally private, satisfying Definition 8.2.

Let $p \in \mathbb{R}_+^k$, $p^T \mathbf{1} = 1$ indicate the true probabilities $p_i = \mathbb{P}(X = i)$. Then by inspection, we have

$$\mathbb{P}(Z = i) = p_i \frac{e^\varepsilon}{k-1+e^\varepsilon} + (1 - p_i) \frac{1}{k-1+e^\varepsilon} = p_i \frac{e^\varepsilon - 1}{e^\varepsilon + k - 1} + \frac{1}{e^\varepsilon + k - 1}.$$

Thus, letting $\hat{c}_n \in \mathbb{R}_+^k$ denote the empirical proportion of the Z observations in a sample of size n , we have

$$\hat{p}_n := \frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \left(\hat{c}_n - \frac{1}{e^\varepsilon + k - 1} \mathbf{1} \right)$$

satisfies $\mathbb{E}[\hat{p}_n] = p$, and we also have

$$\mathbb{E} [\|\hat{p}_n - p\|_2^2] = \left(\frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \right)^2 \mathbb{E} [\|\hat{c}_n - \mathbb{E}[\hat{c}_n]\|_2^2] = \frac{1}{n} \left(\frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \right)^2 \sum_{j=1}^k \mathbb{P}(Z = j)(1 - \mathbb{P}(Z = j)).$$

As $\sum_j \mathbb{P}(Z = j) = 1$, we always have the bound $\mathbb{E}[\|\hat{p}_n - p\|_2^2] \leq \frac{1}{n} \left(\frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \right)^2$.

We may consider two regimes for simplicity: when $\varepsilon \leq 1$ and when $\varepsilon \geq \log k$. In the former case—the high privacy regime—we have $\frac{1}{k} \lesssim \mathbb{P}(Z = i) \lesssim \frac{1}{k}$, so that the mean ℓ_2 squared error scales as $\frac{1}{n} \frac{k^2}{\varepsilon^2}$. When $\varepsilon \geq \log k$ is large, by contrast, we see that the error scales at worst as $\frac{1}{n}$, which is the “non-private” mean squared error. \diamond

While randomized response is essentially the standard mechanism in locally private settings, in centralized privacy, the “standard” mechanism is Laplace noise addition because of its exponential tails. In this case, we require a few additional definitions. Suppose that we wish to release some d -dimensional function $f(P_n)$ of the sample distribution P_n (equivalently, the associated sample X_1^n), where f takes values in \mathbb{R}^d . In the case that f is Lipschitz with respect to the Hamming metric—that is, the counting metric on \mathcal{X}^n —it is relatively straightforward to develop private mechanisms. To better reflect the nomenclature in the privacy literature and easier use in our future development, for $p \in [1, \infty]$ we define the *global sensitivity* of f by

$$\text{GS}_p(f) := \sup_{P_n, P'_n \in \mathcal{P}_n} \left\{ \|f(P_n) - f(P'_n)\|_p \mid d_{\text{ham}}(P_n, P'_n) \leq 1 \right\}.$$

This is simply the Lipschitz constant of f with respect to the Hamming metric. The global sensitivity is a convenient metric, because it allows simple noise addition strategies.

Example 8.1.3 (Laplace mechanisms): Recall the Laplace distribution, parameterized by a shape parameter β , which has density on \mathbb{R} defined by

$$p(w) = \frac{1}{2\beta} \exp(-|w|/\beta),$$

and the analogous d -dimensional variant, which has density

$$p(w) = \frac{1}{(2\beta)^2} \exp(-\|w\|_1 / \beta).$$

If $W \sim \text{Laplace}(\beta)$, $W \in \mathbb{R}$, then $\mathbb{E}[W] = 0$ by symmetry, while $\mathbb{E}[W^2] = \frac{1}{\beta} \int_0^\infty w^2 e^{-w/\beta} = 2\beta^2$.

Suppose that $f : \mathcal{P}_n \rightarrow \mathbb{R}^d$ has finite global sensitivity for the ℓ_1 -norm,

$$\text{GS}_1(f) = \sup \left\{ \|f(P_n) - f(P'_n)\|_1 \mid d_{\text{ham}}(P_n, P'_n) \leq 1, P_n, P'_n \in \mathcal{P}_n \right\}.$$

Letting $L = \text{GS}_1(f)$ be the Lipschitz constant for simplicity, if we consider the mechanism defined by the addition of $W \in \mathbb{R}^d$ with independent $\text{Laplace}(L/\varepsilon)$ coordinates,

$$Z := f(P_n) + W, \quad W_j \stackrel{\text{iid}}{\sim} \text{Laplace}(L/\varepsilon), \quad (8.1.3)$$

we have that Z is ε -differentially private. Indeed, for samples P_n, P'_n differing in at most a single example, Z has density ratio

$$\frac{q(z \mid P_n)}{q(z \mid P'_n)} = \exp \left(-\frac{\varepsilon}{L} \|f(P_n) - z\|_1 + \frac{\varepsilon}{L} \|f(P'_n) - z\|_1 \right) \leq \exp \left(\frac{\varepsilon}{L} \|f(P_n) - f(P'_n)\|_1 \right) \leq \exp(\varepsilon)$$

by the triangle inequality and that f is L -Lipschitz with respect to the Hamming metric. Thus Z is ε -differentially private. Moreover, we have

$$\mathbb{E}[\|Z - f(P_n)\|_2^2] = \frac{2d\text{GS}_1(f)^2}{\varepsilon^2},$$

so that if L is small, we may report the value of f accurately. \diamond

The most common instances and applications of the Laplace mechanism are in estimation of means and histograms. Let us demonstrate more carefully worked examples in these two cases.

Example 8.1.4 (Private one-dimensional mean estimation): Suppose that we have variables X_i taking values in $[-b, b]$ for some $b < \infty$, and wish to estimate $\mathbb{E}[X]$. A natural function to release is then $f(X_1^n) = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. This has Lipschitz constant $2b/n$ with respect to the Hamming metric, because for any two samples $x, x' \in [-b, b]^n$ differing in only entry i , we have

$$|f(x) - f(x')| = \frac{1}{n} |x_i - x'_i| \leq \frac{2b}{n}$$

because $x_i \in [-b, b]$. Thus the Laplace mechanism (8.1.3) with the choice variance $W \sim \text{Laplace}(2b/(n\varepsilon))$ yields

$$\mathbb{E}[(Z - \mathbb{E}[X])^2] = \mathbb{E}[(\bar{X}_n - \mathbb{E}[X])^2] + \mathbb{E}[(Z - \bar{X}_n)^2] = \frac{1}{n} \text{Var}(X) + \frac{8b^2}{n^2 \varepsilon^2} \leq \frac{b^2}{n} + \frac{8b^2}{n^2 \varepsilon^2}.$$

We can privately release means with little penalty so long as $\varepsilon \gg n^{-1/2}$. \diamond

Example 8.1.5 (Private histogram (multinomial) release): Suppose that we wish to estimate a multinomial distribution, or put differently, a histogram. That is, we have observations $X \in \{1, \dots, k\}$, where k may be large, and wish to estimate $p_j := \mathbb{P}(X = j)$ for $j = 1, \dots, k$. For a given sample x_1^n , the empirical count vector \hat{p}_n with coordinates $\hat{p}_{n,j} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i = j\}$ satisfies

$$\text{GS}_1(\hat{p}_n) = \frac{2}{n}$$

because swapping a single example x_i for x'_i may change the counts for at most two coordinates j, j' by 1. Consequently, the Laplace noise addition mechanism

$$Z = \hat{p}_n + W, \quad W_j \stackrel{\text{iid}}{\sim} \text{Laplace}\left(\frac{2}{n\varepsilon}\right)$$

satisfies

$$\mathbb{E}[\|Z - \hat{p}_n\|_2^2] = \frac{8k}{n^2 \varepsilon^2}$$

and consequently

$$\mathbb{E}[\|Z - p\|_2^2] = \frac{8k}{n^2 \varepsilon^2} + \frac{1}{n} \sum_{j=1}^k p_j(1 - p_j) \leq \frac{8k}{n^2 \varepsilon^2} + \frac{1}{n}.$$

This example shows one of the challenges of differentially private mechanisms: even in the case where the quantity of interest is quite stable (insensitive to changes in the underlying sample, or has small Lipschitz constant), it may be the case that the resulting mechanism adds noise that introduces some dimension-dependent scaling. In this case, the conditions on privacy levels acceptable for good estimation—in that the rate of convergence is no different from the non-private case, which achieves $\mathbb{E}[\|\hat{p}_n - p\|_2^2] = \frac{1}{n} \sum_{j=1}^k p_j(1 - p_j) \leq \frac{1}{n}$ are that $\varepsilon \gg \frac{k}{n}$. Thus, in the case that the histogram has a large number of bins, the naive noise addition strategy cannot provide as much protection without sacrificing efficiency.

If instead of ℓ_2 -error we consider ℓ_∞ error, it is possible to provide somewhat more satisfying results in this case. Indeed, we know that $\mathbb{P}(\|W\|_\infty \geq t) \leq k \exp(-t/b)$ for $W_j \stackrel{\text{iid}}{\sim} \text{Laplace}(b)$, so that in the mechanism above we have

$$\mathbb{P}(\|Z - \hat{p}_n\|_\infty \geq t) \leq k \exp\left(-\frac{tn\varepsilon}{2}\right) \quad \text{all } t \geq 0,$$

so using that each coordinate of \hat{p}_n is 1-sub-Gaussian, we have

$$\begin{aligned}\mathbb{E}[\|Z - p\|_\infty] &\leq \mathbb{E}[\|\hat{p}_n - p\|_\infty] + \mathbb{E}[\|W\|_\infty] \leq \sqrt{\frac{2 \log k}{n}} + \inf_{t \geq 0} \left\{ t + \frac{2k}{n\varepsilon} \exp\left(-\frac{tn\varepsilon}{2}\right) \right\} \\ &\leq \sqrt{\frac{2 \log k}{n}} + \frac{2 \log k}{n\varepsilon} + \frac{2}{n\varepsilon}.\end{aligned}$$

In this case, then, whenever $\varepsilon \gg (n/\log k)^{-1/2}$, we obtain rate of convergence at least $\sqrt{2 \log k/n}$, which is a bit loose (as we have not controlled the variance of \hat{p}_n), but somewhat more satisfying than the k -dependent penalty above. \diamond

8.1.2 Resilience to side information, Bayesian perspectives, and data processing

One of the major challenges in the definition of privacy is to protect against side information, especially because in the future, information about you may be compromised, allowing various linkage attacks. With this in mind, we return to our three desiderata. First, we note the following simple fact: if Z is a differentially private view of a sample X_1^n (or associated empirical distribution P_n), then any downstream functions Y are also differentially private. That is, if we have the Markov chain $P_n \rightarrow Z \rightarrow Y$, then for any $P_n, P'_n \in \mathcal{P}_n$ with $d_{\text{ham}}(P_n, P'_n) \leq 1$, we have for any set A that

$$\frac{\mathbb{P}(Y \in A \mid x)}{\mathbb{P}(Y \in A \mid x')} = \frac{\int P(Y \in A \mid z) q(z \mid P_n) d\mu(z)}{\int P(Y \in A \mid z) q(z \mid P'_n) d\mu(z)} \leq e^\varepsilon \frac{\int P(Y \in A \mid z) q(z \mid P'_n) d\mu(z)}{\int P(Y \in A \mid z) q(z \mid P'_n) d\mu(z)} = e^\varepsilon.$$

That is, any type of post-processing cannot reduce privacy.

With this simple idea out of the way, let us focus on our testing-based desideratum. In this case, we consider a testing scenario, where an adversary wishes to test two hypotheses against one another, where the hypotheses are

$$H_0 : X_1^n = x_1^n \quad \text{vs.} \quad H_1 : X_1^n = (x_1^{i-1}, x'_i, x_{i+1}^n),$$

so that the samples under H_0 and H_1 differ only in the i th observation $X_i \in \{x_i, x'_i\}$. Now, for a channel taking inputs from \mathcal{X}^n and outputting $Z \in \mathcal{Z}$, we define ε -conditional hypothesis testing privacy by saying that

$$Q(\Psi(Z) = 1 \mid H_0, Z \in A) + Q(\Psi(Z) = 0 \mid H_1, Z \in A) \geq 1 - \varepsilon \quad (8.1.4)$$

for all sets $A \subset \mathcal{Z}$ satisfying $Q(A \mid H_0) > 0$ and $Q(A \mid H_1) > 0$. That is, roughly, no matter *what* value Z takes on, the probability of error in a test of whether H_0 or H_1 is true—even with knowledge of $x_j, j \neq i$ —is high. We then have the following proposition.

Proposition 8.1.6. *Assume the channel Q is ε -differentially private. Then Q is also $\bar{\varepsilon} = 1 - e^{-2\varepsilon} \leq 2\varepsilon$ -conditional hypothesis testing private.*

Proof Let Ψ be any test of H_0 versus H_1 , and let $B = \{z \mid \Psi(z) = 1\}$ be the acceptance region of the test. Then

$$\begin{aligned}Q(B \mid H_0, Z \in A) + Q(B^c \mid H_1, Z \in A) &= \frac{Q(A, B \mid H_0)}{Q(A \mid H_0)} + \frac{Q(A, B^c \mid H_1)}{Q(A \mid H_1)} \\ &\geq e^{-2\varepsilon} \frac{Q(A, B \mid H_1)}{Q(A \mid H_1)} + \frac{Q(A, B^c \mid H_1)}{Q(A \mid H_1)} \\ &\geq e^{-2\varepsilon} \frac{Q(A, B \mid H_1) + Q(A, B^c \mid H_1)}{Q(A \mid H_1)},\end{aligned}$$

where the first inequality uses ε -differential privacy. Then we simply note that $Q(A, B \mid H_1) + Q(A, B^c \mid H_1) = Q(A \mid H_1)$. \square

So we see that (roughly), even conditional on the output of the channel, we still cannot test whether the initial dataset was x or x' whenever x, x' differ in only a single observation.

An alternative perspective is to consider a Bayesian one, which allows us to more carefully consider side information. In this case, we consider the following thought experiment. An adversary has a set of prior beliefs π on \mathcal{X}^n , and we consider the adversary's posterior $\pi(\cdot \mid Z)$ induced by observing the output Z of some mechanism M . In this case, *Bayes factors*, which measure how much prior and posterior distributions differ after observations, provide one immediate perspective.

Proposition 8.1.7. *A mechanism $M : \mathcal{P}_n \rightarrow \mathcal{Z}$ is ε -differentially private if and only if for any prior distribution π on \mathcal{P}_n and any observation $z \in \mathcal{Z}$, the posterior odds satisfy*

$$\frac{\pi(P_n \mid z)}{\pi(P'_n \mid z)} \leq e^\varepsilon$$

for all $P_n, P'_n \in \mathcal{P}_n$ with $d_{\text{ham}}(P_n, P'_n) \leq 1$.

Proof Let q be the associated density of $Z = M(\cdot)$ (conditional or marginal). We have $\pi(P_n \mid z) = q(z \mid P_n)\pi(P_n)/q(z)$. Then

$$\frac{\pi(P_n \mid z)}{\pi(P'_n \mid z)} = \frac{q(z \mid P_n)\pi(P_n)}{q(z \mid P'_n)\pi(P'_n)} \leq e^\varepsilon \frac{\pi(P_n)}{\pi(P'_n)}$$

for all z, P_n, P'_n if and only if M is ε -differentially private. \square

Thus we see that private channels mean that prior and posterior odds between two neighboring samples cannot change substantially, no matter what the observation Z actually is.

For an alternative view, we consider a somewhat restricted family of prior distributions, where we now take the view of a sample $x_1^n \in \mathcal{X}^n$. There is some annoyance in this calculation in that the *order* of the sample may be important, but it at least gets toward some semantic interpretation of differential privacy. We consider the adversary's beliefs on whether a particular value x belongs to the sample, but more precisely, we consider whether $X_i = x$. We assume that the prior density π on \mathcal{X}^n satisfies

$$\pi(x_1^n) = \pi_{\setminus i}(x_{\setminus i})\pi_i(x_i), \quad (8.1.5)$$

where $x_{\setminus i} = (x_1^{i-1}, x_{i+1}^n) \in \mathcal{X}^{n-1}$. That is, the adversary's beliefs about person i in the dataset are independent of his beliefs about the other members of the dataset. (We assume that π is a density with respect to a measure μ on $\mathcal{X}^{n-1} \times \mathcal{X}$, where $d\mu(s, x) = d\mu(s)d\mu(x)$.) Under the condition (8.1.5), we have the following proposition.

Proposition 8.1.8. *Let Q be an ε -differentially private channel and let π be any prior distribution satisfying condition (8.1.5). Then for any z , the posterior density π_i on X_i satisfies*

$$e^{-\varepsilon}\pi_i(x) \leq \pi_i(x \mid Z = z) \leq e^\varepsilon\pi_i(x).$$

Proof We abuse notation and for a sample $s \in \mathcal{X}^{n-1}$, where $s = (x_1^{i-1}, x_{i+1}^n)$, we let $s \oplus_i x = (x_1^{i-1}, x, x_{i+1}^n)$. Letting μ be the base measure on $\mathcal{X}^{n-1} \times \mathcal{X}$ with respect to which π is a density and $q(\cdot | x_1^n)$ be the density of the channel Q , we have

$$\begin{aligned} \pi_i(x | Z = z) &= \frac{\int_{s \in \mathcal{X}^{n-1}} q(z | s \oplus_i x) \pi(s \oplus_i x) d\mu(s)}{\int_{s \in \mathcal{X}^{n-1}} \int_{x' \in \mathcal{X}} q(z | s \oplus_i x') \pi(s \oplus_i x') d\mu(s, x')} \\ &\stackrel{(\star)}{\leq} e^\varepsilon \frac{\int_{s \in \mathcal{X}^{n-1}} q(z | s \oplus_i x) \pi(s \oplus_i x) d\mu(s)}{\int_{s \in \mathcal{X}^{n-1}} \int_{x' \in \mathcal{X}} q(z | s \oplus_i x') \pi(s \oplus_i x') d\mu(s) d\mu(x')} \\ &= e^\varepsilon \frac{\int_{s \in \mathcal{X}^{n-1}} q(z | s \oplus_i x) \pi_{\setminus i}(s) d\mu(s) \pi_i(x)}{\int_{s \in \mathcal{X}^{n-1}} q(z | s \oplus_i x) \pi_{\setminus i}(s) d\mu(s) \int_{x' \in \mathcal{X}} \pi_i(x') d\mu(x')} \\ &= e^\varepsilon \pi_i(x), \end{aligned}$$

where inequality (\star) follows from ε -differential privacy. The lower bound is similar. \square

Roughly, however, we see that Proposition 8.1.8 captures the idea that even if an adversary has substantial prior knowledge—in the form of a prior distribution π on the i th value X_i and everything else in the sample—the posterior cannot change much.

8.2 Weakenings of differential privacy

One challenge with the definition of differential privacy is that it can sometimes require the addition of more noise to a desired statistic than is practical for real use. Moreover, the privacy considerations interact in different ways with geometry: as we saw in Example 8.1.5, the Laplace mechanism adds noise that introduces dimension-dependent scaling, which we discuss more in Example 8.2.9. Consequently, it is of interest to develop weaker notions that—at least hopefully—still provide appropriate and satisfactory privacy protections. To that end, we develop two additional types of privacy that allow the development of more sophisticated and lower-noise mechanisms than standard differential privacy; their protections are necessarily somewhat weaker but are typically satisfactory.

We begin with a definition that allows (very rare) catastrophic privacy breaches—as long as the probability of this event is extremely small (say, 10^{-20}), these may be acceptable.

Definition 8.3. Let $\varepsilon, \delta \geq 0$. A mechanism $M : \mathcal{P}_n \rightarrow \mathcal{Z}$ is (ε, δ) -differentially private if for all (measurable) sets $S \subset \mathcal{Z}$ and all neighboring samples P_n, P'_n ,

$$\mathbb{P}(M(P_n) \in S) \leq e^\varepsilon \mathbb{P}(M(P'_n) \in S) + \delta. \quad (8.2.1)$$

One typically thinks of δ in the definition above as satisfying $\delta = \delta_n$, where $\delta_n \ll n^{-k}$ for any $k \in \mathbb{N}$. (That is, δ decays super-polynomially to zero.) Some practitioners contend that all real-world differentially private algorithms are in fact (ε, δ) -differentially private: while one may use cryptographically secure random number generators, there is some possibility (call this δ) that a cryptographic key may leak, or an encoding may be broken, in the future, making any mechanism (ε, δ) -private at best for some $\delta > 0$.

An alternative definition of privacy is based on Rényi divergences between distributions. These are essentially simply monotonically transformed f divergences (recall Chapter 2.2), though their structure is somewhat more amenable to analysis, especially in our contexts. With that in mind, we define

Definition 8.4. Let P and Q be distributions on a space \mathcal{X} with densities p and q (with respect to a measure μ). For $\alpha \in [1, \infty]$, the Rényi- α -divergence between P and Q is

$$D_\alpha(P\|Q) := \frac{1}{\alpha - 1} \log \int \left(\frac{p(x)}{q(x)} \right)^\alpha q(x) d\mu(x).$$

Here, the values $\alpha \in \{1, \infty\}$ are defined in terms of their respective limits.

Rényi divergences satisfy $\exp((\alpha - 1)D_\alpha(P\|Q)) = 1 + D_f(P\|Q)$, i.e., $D_\alpha(P\|Q) = \frac{1}{\alpha - 1} \log(1 + D_f(P\|Q))$, for the f -divergence defined by $f(t) = t^\alpha - 1$, so that they inherit a number of the properties of such divergences. We enumerate a few here for later reference.

Proposition 8.2.1 (Basic facts on Rényi divergence). *Rényi divergences satisfy the following.*

- i. The divergence $D_\alpha(P\|Q)$ is non-decreasing in α .
- ii. $\lim_{\alpha \downarrow 1} D_\alpha(P\|Q) = D_{\text{kl}}(P\|Q)$ and $\lim_{\alpha \uparrow \infty} D_\alpha(P\|Q) = \sup\{t \mid Q(p(X)/q(X) \geq t) > 0\}$.
- iii. Let $K(\cdot \mid x)$ be a Markov kernel from $\mathcal{X} \rightarrow \mathcal{Z}$ as in Proposition 2.2.13, and let K_P and K_Q be the induced marginals of P and Q under K , respectively. Then $D_\alpha(K_P\|K_Q) \leq D_\alpha(P\|Q)$.

We leave the proof of this proposition as Exercise 8.1, noting that property i is a consequence of Hölder's inequality, property ii is by L'Hopital's rule, and property iii is an immediate consequence of Proposition 2.2.13. Rényi divergences also tensorize nicely—generalizing the tensorization properties of KL-divergence and information of Chapter 2 (recall the chain rule (2.1.6) for KL-divergence)—and we return to this later. As a preview, however, these tensorization properties allow us to prove that the composition of multiple private data releases remains appropriately private.

With these preliminaries in place, we can then provide

Definition 8.5 (Rényi-differential privacy). Let $\varepsilon \geq 0$ and $\alpha \in [1, \infty]$. A channel Q from \mathcal{P}_n to output space \mathcal{Z} is (ε, α) -Rényi private if for all neighboring samples $P_n, P'_n \in \mathcal{P}_n$,

$$D_\alpha(Q(\cdot \mid P_n)\|Q(\cdot \mid P'_n)) \leq \varepsilon. \quad (8.2.2)$$

Clearly, any ε -differentially private channel is also (ε, α) -Rényi private for any $\alpha \geq 1$; as we soon see, we can provide tighter guarantees than this.

8.2.1 Basic mechanisms

We now describe a few of the basic mechanisms that provide guarantees of (ε, δ) -differential privacy and (ε, α) -Rényi privacy. The advantage for these settings is that they allow mechanisms that more naturally handle vectors in ℓ_2 , and smoothness with respect to Euclidean norms, than with respect to ℓ_1 , which is most natural for pure ε -differential privacy. A starting point is the following example, which we will leverage frequently.

Example 8.2.2 (Rényi divergence between Gaussian distributions): Consider normal distributions $\mathcal{N}(\mu_0, \Sigma)$ and $\mathcal{N}(\mu_1, \Sigma)$. Then

$$D_\alpha(\mathcal{N}(\mu_0, \Sigma)\|\mathcal{N}(\mu_1, \Sigma)) = \frac{\alpha}{2}(\mu_0 - \mu_1)^T \Sigma^{-1}(\mu_0 - \mu_1). \quad (8.2.3)$$

To see this equality, we compute the appropriate integral of the densities. Let p and q be the densities of $\mathbf{N}(\mu_0, \Sigma)$ and $\mathbf{N}(\mu_1, \Sigma)$, respectively. Then letting \mathbb{E}_{μ_1} denote expectation over $X \sim \mathbf{N}(\mu_1, \Sigma)$, we have

$$\begin{aligned} \int \left(\frac{p(x)}{q(x)} \right)^\alpha q(x) dx &= \mathbb{E}_{\mu_1} \left[\exp \left(-\frac{\alpha}{2} (X - \mu_0)^T \Sigma^{-1} (X - \mu_0) + \frac{\alpha}{2} (X - \mu_1)^T \Sigma^{-1} (X - \mu_1) \right) \right] \\ &\stackrel{(i)}{=} \mathbb{E}_{\mu_1} \left[\exp \left(-\frac{\alpha}{2} (\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1) + \alpha (\mu_0 - \mu_1)^T \Sigma^{-1} (X - \mu_1) \right) \right] \\ &\stackrel{(ii)}{=} \exp \left(-\frac{\alpha}{2} (\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1) + \frac{\alpha^2}{2} (\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1) \right), \end{aligned}$$

where equality (i) is simply using that $(x - a)^2 - (x - b)^2 = (a - b)^2 + 2(b - a)(x - b)$ and equality (ii) follows because $(\mu_0 - \mu_1)^T \Sigma^{-1} (X - \mu_1) \sim \mathbf{N}(0, (\mu_1 - \mu_0)^T \Sigma^{-1} (\mu_1 - \mu_0))$ under $X \sim \mathbf{N}(\mu_1, \Sigma)$. Noting that $-\alpha + \alpha^2 = \alpha(\alpha - 1)$ and taking logarithms gives the result. \diamond

Example 8.2.2 is the key to developing different privacy-preserving schemes under Rényi privacy. Let us reconsider Example 8.1.3, except that instead of assuming the function f of interest is smooth with respect to ℓ_1 norm, we use the ℓ_2 -norm.

Example 8.2.3 (Gaussian mechanisms): Suppose that $f : \mathcal{P}_n \rightarrow \mathbb{R}^d$ has Lipschitz constant L with respect to the ℓ_2 -norm (for the Hamming metric d_{ham}), that is, global ℓ_2 -sensitivity

$$\text{GS}_2(f) = \sup \{ \|f(P_n) - f(P'_n)\|_2 \mid d_{\text{ham}}(P_n, P'_n) \leq 1 \} \leq L.$$

Then, for any variance $\sigma^2 > 0$, we have that the mechanism

$$Z = f(P_n) + W, \quad W \sim \mathbf{N}(0, \sigma^2 I)$$

satisfies

$$D_\alpha(\mathbf{N}(f(P_n), \sigma^2) \parallel \mathbf{N}(f(P'_n), \sigma^2)) = \frac{\alpha}{2\sigma^2} \|f(P_n) - f(P'_n)\|_2^2 \leq \frac{\alpha}{2\sigma^2} L^2$$

for neighboring samples P_n, P'_n . Thus, if we have Lipschitz constant L and desire (ε, α) -Rényi privacy, we may take $\sigma^2 = \frac{L^2 \alpha}{2\varepsilon}$, and then the mechanism

$$Z = f(P_n) + W \quad W \sim \mathbf{N}\left(0, \frac{L^2 \alpha}{2\varepsilon} I\right) \tag{8.2.4}$$

satisfies (ε, α) -Rényi privacy. \diamond

Certain special cases can make this more concrete. Indeed, suppose we wish to estimate a mean $\mathbb{E}[X]$ where $X_i \stackrel{\text{iid}}{\sim} P$ for some distribution P such that $\|X_i\|_2 \leq r$ with probability 1 for some radius.

Example 8.2.4 (Bounded mean estimation with Gaussian mechanisms): Letting $f(X_1^n) = \bar{X}_n$ be the sample mean, where X_i satisfy $\|X_i\|_2 \leq r$ as above, we see immediately that

$$\text{GS}_2(f) = \frac{2r}{n}.$$

In this case, the Gaussian mechanism (8.2.4) with $L = \frac{2r}{n}$ yields

$$\mathbb{E} \left[\|Z - \bar{X}_n\|_2^2 \right] = \mathbb{E}[\|W\|_2^2] = \frac{2dr^2\alpha}{n^2\varepsilon}.$$

Then we have

$$\mathbb{E}[\|Z - \mathbb{E}[X]\|_2^2] = \mathbb{E}[\|\bar{X}_n - \mathbb{E}[X]\|_2^2] + \mathbb{E}[\|Z - \bar{X}_n\|_2^2] \leq \frac{r^2}{n} + \frac{2dr^2\alpha}{n^2\varepsilon}.$$

It is not immediately apparent how to compare this quantity to the case for the Laplace mechanism in Example 8.1.3, but we will return to this shortly once we have developed connections between the various privacy notions we have developed. \diamond

8.2.2 Connections between privacy measures

An important consideration in our development of privacy definitions and mechanisms is to understand the relationships between the definitions, and when a channel Q satisfying one of the definitions satisfies one of our other definitions. Thus, we collect a few different consequences of our definitions, which help to show the various definitions are stronger or weaker than others.

First, we argue that ε -differential privacy implies stronger values of Rényi-differential privacy.

Proposition 8.2.5. *Let $\varepsilon \geq 0$ and let P and Q be distributions such that $e^{-\varepsilon} \leq P(A)/Q(A) \leq e^\varepsilon$ for all measurable sets A . Then for any $\alpha \in [1, \infty]$,*

$$D_\alpha(P\|Q) \leq \min \left\{ \frac{3\alpha}{2}\varepsilon^2, \varepsilon \right\}.$$

As an immediate corollary, we have

Corollary 8.2.6. *Let $\varepsilon \geq 0$ and assume that Q is ε -differentially private. Then for any $\alpha \geq 1$, Q is $(\min\{\frac{3\alpha}{2}\varepsilon^2, \varepsilon\}, \alpha)$ -Rényi private.*

Before proving the proposition, let us see its implications for Example 8.2.4 versus estimation under ε -differential privacy. Let $\varepsilon \leq 1$, so that roughly to have “similar” privacy, we require that our Rényi private channels satisfy $D_\alpha(Q(\cdot | x)\|Q(\cdot | x')) \leq \varepsilon^2$. The ℓ_1 -sensitivity of the mean satisfies $\|\bar{x}_n - \bar{x}'_n\|_1 \leq \sqrt{d}\|\bar{x}_n - \bar{x}'_n\|_2 \leq 2\sqrt{d}r/n$ for neighboring samples. Then the Laplace mechanism (8.1.3) satisfies

$$\mathbb{E}[\|Z_{\text{Laplace}} - \mathbb{E}[X]\|_2^2] = \mathbb{E}[\|\bar{X}_n - \mathbb{E}[X]\|_2^2] + \frac{8r^2}{n^2\varepsilon^2} \cdot d^2,$$

while the Gaussian mechanism under (ε^2, α) -Rényi privacy will yield

$$\mathbb{E}[\|Z_{\text{Gauss}} - \mathbb{E}[X]\|_2^2] = \mathbb{E}[\|\bar{X}_n - \mathbb{E}[X]\|_2^2] + \frac{2r^2}{n^2\varepsilon^2} \cdot d\alpha.$$

This is evidently better than the Laplace mechanism whenever $\alpha < d$.

Proof of Proposition 8.2.5 We assume that P and Q have densities p and q with respect to a base measure μ , which is no loss of generality, whence the ratio condition implies that $e^{-\varepsilon} \leq p/q \leq e^\varepsilon$ and $D_\alpha(P\|Q) = \frac{1}{\alpha-1} \log \int (p/q)^\alpha q d\mu$. We prove the result assuming that $\alpha \in (1, \infty)$, as continuity gives the result for $\alpha \in \{1, \infty\}$.

First, it is clear that $D_\alpha(P\|Q) \leq \varepsilon$ always. For the other term in the minimum, let us assume that $\alpha \leq 1 + \frac{1}{\varepsilon}$ and $\varepsilon \leq 1$. If either of these fails, the result is trivial, because for $\alpha > 1 + \frac{1}{\varepsilon}$ we have $\frac{3}{2}\alpha\varepsilon^2 \geq \frac{3}{2}\varepsilon \geq \varepsilon$, and similarly $\varepsilon \geq 1$ implies $\frac{3}{2}\alpha\varepsilon^2 \geq \varepsilon$.

Now we perform a Taylor approximation of $t \mapsto (1+t)^\alpha$. By Taylor's theorem, we have for any $t > -1$ that

$$(1+t)^\alpha = 1 + \alpha t + \frac{\alpha(\alpha-1)}{2}(1+\tilde{t})^{\alpha-2}t^2$$

for some $\tilde{t} \in [0, t]$ (or $[t, 0]$ if $t < 0$). In particular, if $1+t \leq c$, then $(1+t)^\alpha \leq 1 + \alpha t + \frac{\alpha(\alpha-1)}{2} \max\{1, c^{\alpha-2}\}t^2$. Now, we compute the divergence: we have

$$\begin{aligned} \exp((\alpha-1)D_\alpha(P\|Q)) &= \int \left(\frac{p(z)}{q(z)}\right)^\alpha q(z) d\mu(z) \\ &= \int \left(1 + \frac{p(z)}{q(z)} - 1\right)^\alpha q(z) d\mu(z) \\ &\leq 1 + \alpha \int \left(\frac{p(z)}{q(z)} - 1\right) q(z) d\mu(z) + \frac{\alpha(\alpha-1)}{2} \max\{1, \exp(\varepsilon(\alpha-2))\} \int \left(\frac{p(z)}{q(z)} - 1\right)^2 q(z) d\mu(z) \\ &\leq 1 + \frac{\alpha(\alpha-1)}{2} e^{\varepsilon[\alpha-2]_+} \cdot (e^\varepsilon - 1)^2. \end{aligned}$$

Now, we know that $\alpha - 2 \leq 1/\varepsilon - 1$ by assumption, so using that $\log(1+x) \leq x$, we obtain

$$D_\alpha(P\|Q) \leq \frac{\alpha}{2} (e^\varepsilon - 1)^2 \cdot \exp([1 - \varepsilon]_+).$$

Finally, a numerical calculation yields that this quantity is at most $\frac{3\alpha}{2}\varepsilon^2$ for $\varepsilon \leq 1$. \square

We can also provide connections from (ε, α) -Rényi privacy to (ε, δ) -differential privacy, and then from there to ε -differential privacy. We begin by showing how to develop (ε, δ) -differential privacy out of Rényi privacy. Another way to think about this proposition is that whenever two distributions P and Q are close in Rényi divergence, then there is some limited “amplification” of probabilities that is possible in moving from one to the other.

Proposition 8.2.7. *Let P and Q satisfy $D_\alpha(P\|Q) \leq \varepsilon$. Then for any set A ,*

$$P(A) \leq \exp\left(\frac{\alpha-1}{\alpha}\varepsilon\right) Q(A)^{\frac{\alpha-1}{\alpha}}.$$

Consequently, for any $\delta > 0$,

$$P(A) \leq \min\left\{\exp\left(\varepsilon + \frac{1}{\alpha-1} \log \frac{1}{\delta}\right) Q(A), \delta\right\} \leq \exp\left(\varepsilon + \frac{1}{\alpha-1} \log \frac{1}{\delta}\right) Q(A) + \delta.$$

As above, we have an immediate corollary to this result.

Corollary 8.2.8. *Assume that M is (ε, α) -Rényi private. Then it is also $(\varepsilon + \frac{1}{\alpha-1} \log \frac{1}{\delta}, \delta)$ -differentially private for any $\delta > 0$.*

Before turning to the proof of the proposition, we show how it can provide prototypical (ε, δ) -private mechanisms via Gaussian noise addition.

Example 8.2.9 (Gaussian mechanisms, continued): Consider Example 8.2.3, where $f : \mathcal{P}_n \rightarrow \mathbb{R}^d$ has ℓ_2 -sensitivity L . Then by Example 8.2.2, the Gaussian mechanism $Z = f(P_n) + W$ for $W \sim \mathcal{N}(0, \sigma^2 I)$ is $(\frac{\alpha L^2}{2\sigma^2}, \alpha)$ -Rényi private for all $\alpha \geq 1$. Combining this with Corollary 8.2.8, the Gaussian mechanism is also

$$\left(\frac{\alpha L^2}{2\sigma^2} + \frac{1}{\alpha - 1} \log \frac{1}{\delta}, \delta \right)\text{-differentially private}$$

for any $\delta > 0$ and $\alpha > 1$. Optimizing first over α by taking $\alpha = 1 + \sqrt{2\sigma^2 \log \delta^{-1}/L^2}$, we see that the channel is $(\frac{L^2}{2\sigma^2} + \sqrt{2L^2 \log \delta^{-1}/\sigma^2}, \delta)$ -differentially private. Thus we have that the Gaussian mechanism

$$Z = f(P_n) + W, \quad W \sim \mathcal{N}(0, \sigma^2 I) \text{ for } \sigma^2 = L^2 \max \left\{ \frac{8 \log \frac{1}{\delta}}{\varepsilon^2}, \frac{1}{\varepsilon} \right\} \quad (8.2.5)$$

is (ε, δ) -differentially private.

To continue with our ℓ_2 -bounded mean-estimation in Example 8.2.4, let us assume that $\varepsilon < 8 \log \frac{1}{\delta}$, in which case the Gaussian mechanism (8.2.5) with $L^2 = r^2/n^2$ achieves (ε, δ) -differential privacy, and we have

$$\mathbb{E}[\|Z_{\text{Gauss}} - \mathbb{E}[X]\|_2^2] = \mathbb{E}[\|\bar{X}_n - \mathbb{E}[X]\|_2^2] + O(1) \frac{r^2}{n^2 \varepsilon^2} \cdot d \log \frac{1}{\delta}.$$

Comparing to the previous cases, we see an improvement over the Laplace mechanism whenever $\log \frac{1}{\delta} \ll d$, or that $\delta \gg e^{-d}$. \diamond

Proof of Proposition 8.2.7 We use the data processing inequality of Proposition 8.2.1.iii, which shows that

$$\varepsilon \geq D_\alpha(P \| Q) \geq \frac{1}{\alpha - 1} \log \left[\left(\frac{P(A)}{Q(A)} \right)^\alpha Q(A) \right].$$

Rearranging and taking exponentials, we immediately obtain the first claim of the proposition.

For the second, we require a bit more work. First, let us assume that $Q(A) > e^{-\varepsilon} \delta^{\frac{\alpha}{\alpha-1}}$. Then we have by the first claim of the proposition that

$$\begin{aligned} P(A) &\leq \exp \left(\frac{\alpha - 1}{\alpha} \varepsilon + \frac{1}{\alpha} \log \frac{1}{Q(A)} \right) Q(A) \\ &\leq \exp \left(\frac{\alpha - 1}{\alpha} \varepsilon + \frac{1}{\alpha} \varepsilon + \frac{1}{\alpha - 1} \log \frac{1}{\delta} \right) Q(A) = \exp \left(\varepsilon + \frac{1}{\alpha - 1} \log \frac{1}{\delta} \right) Q(A). \end{aligned}$$

On the other hand, when $Q(A) \leq e^{-\varepsilon} \delta^{\frac{\alpha}{\alpha-1}}$, then again using the first result of the proposition,

$$\begin{aligned} P(A) &\leq \exp \left(\frac{\alpha - 1}{\alpha} (\varepsilon + \log Q(A)) \right) \\ &\leq \exp \left(\frac{\alpha - 1}{\alpha} \left(\varepsilon - \varepsilon + \frac{\alpha}{\alpha - 1} \log \delta \right) \right) = \delta. \end{aligned}$$

This gives the second claim of the proposition. \square

Finally, we develop our last set of connections, which show how we may relate (ε, δ) -private channels with ε -private channels. To provide this definition, we require one additional weakened notion of divergence, which relates (ε, δ) -differential privacy to Rényi- α -divergence with $\alpha = \infty$. We define

$$D_{\infty}^{\delta}(P\|Q) := \sup_{S \subset \mathcal{X}} \left\{ \log \frac{P(S) - \delta}{Q(S)} \mid P(S) > \delta \right\},$$

where the supremum is over measurable sets. Evidently equivalent to this definition is that $D_{\infty}^{\delta}(P\|Q) \leq \varepsilon$ if and only if

$$P(S) \leq e^{\varepsilon} Q(S) + \delta \quad \text{for all } S \subset \mathcal{X}.$$

Then we have the following lemma.

Lemma 8.2.10. *Let $\varepsilon > 0$ and $\delta \in (0, 1)$, and let P and Q be distributions on a space \mathcal{X} .*

- (i) *We have $D_{\infty}^{\delta}(P\|Q) \leq \varepsilon$ if and only if there exists a probability distribution R on \mathcal{X} such that $\|P - R\|_{\text{TV}} \leq \delta$ and $D_{\infty}(R\|Q) \leq \varepsilon$.*
- (ii) *We have $D_{\infty}^{\delta}(P\|Q) \leq \varepsilon$ and $D_{\infty}^{\delta}(Q\|P) \leq \varepsilon$ if and only if there exist distributions P_0 and Q_0 such that*

$$\|P - P_0\|_{\text{TV}} \leq \frac{\delta}{1 + e^{\varepsilon}}, \quad \|Q - Q_0\|_{\text{TV}} \leq \frac{\delta}{1 + e^{\varepsilon}},$$

and

$$D_{\infty}(P_0\|Q_0) \leq \varepsilon \quad \text{and} \quad D_{\infty}(Q_0\|P_0) \leq \varepsilon.$$

The proof of the lemma is technical, so we defer it to Section 8.5.1. The key application of the lemma—which we shall see presently—is that (ε, δ) -differentially private algorithms compose in elegant ways.

8.2.3 Side information protections under weakened notions of privacy

We briefly discuss the side information protections these weaker notions of privacy protect. For both (ε, δ) -differential privacy and (ε, α) -Rényi privacy, we revisit the treatment in Proposition 8.1.7, considering Bayes factors and ratios of prior and posterior divergences, as these are natural formulations of side information in terms of an adversary's probabilistic beliefs. Our first analogue of Proposition 8.1.7, applies to the (ε, δ) -private case.

Proposition 8.2.11. *Let M be a (ε, δ) -differentially private mechanism. Then for any neighboring $P_n, P'_n, P_n^{(0)} \in \mathcal{P}_n$, we have with probability at least $1 - \delta$ over the draw of $Z = M(P_n^{(0)})$, the posterior odds satisfy*

$$\frac{\pi(P_n \mid z)}{\pi(P'_n \mid z)} \leq e^{3\varepsilon} \frac{\pi(P_n)}{\pi(P'_n)}.$$

Deferring the proof momentarily, this result shows that as long as two samples x, x' are neighboring, then an adversary is extremely unlikely to be able to glean substantially distinguishing information between the samples. This is suggestive of a heuristic in differential privacy that if n is the sample size, then one should take $\delta \ll 1/n$ to limit the probability of disclosure: by a union bound, we see that for each individual $i \in \{1, \dots, n\}$, we can simultaneously guarantee that the posterior odds for swapping individual i 's data do not change much (with high probability).

Unsurprisingly at this point, we can also give posterior update bounds for Rényi differential privacy. Here, instead of giving high-probability bounds—though it is possible—we can show that moments of the odds ratio do not change significantly. Indeed, we have the following proposition:

Proposition 8.2.12. *Let M be a (ε, α) -Rényi private mechanism, where $\alpha \in (1, \infty)$. Then for any neighboring $P_n, P'_n, P_n^{(0)} \in \mathcal{P}_n$, we have*

$$\mathbb{E}_0 \left[\left(\frac{\pi(P_n | Z)}{\pi(P'_n | Z)} \right)^{\alpha-1} \right]^{\frac{1}{\alpha-1}} \leq e^\varepsilon \frac{\pi(P_n)}{\pi(P'_n)},$$

where \mathbb{E}_0 denotes expectation taken over $Z = M(P_n^{(0)})$.

Proposition 8.2.12 communicates a similar message to our previous results in this vein: even if we get information from the output of the private mechanism on some sample $x_0 \in \mathcal{X}^n$ near the samples (datasets) of interest x, x' that an adversary wishes to distinguish, it is impossible to update beliefs by much. The parameter α then controls the degree of difficulty of this “impossible” claim, which one can see by (for example) applying a Chebyshev-type bound to the posterior ratio and prior ratios.

We now turn to the promised proofs of Propositions 8.2.11 and 8.2.12. To prove the former, we require a definition.

Definition 8.6. *Distributions P and Q on a space \mathcal{X} are (ε, δ) -close if for all measurable A*

$$P(A) \leq e^\varepsilon Q(A) + \delta \quad \text{and} \quad Q(A) \leq e^\varepsilon P(A) + \delta.$$

Letting p and q denote their densities (with respect to any shared base measure), they are (ε, δ) -pointwise close if the set

$$A := \{x \in \mathcal{X} : e^{-\varepsilon} q(x) \leq p(x) \leq e^\varepsilon q(x)\} = \{x \in \mathcal{X} : e^{-\varepsilon} p(x) \leq q(x) \leq e^\varepsilon p(x)\}$$

satisfies $P(A) \geq 1 - \delta$ and $Q(A) \geq 1 - \delta$.

The following lemma shows the strong relationship between closeness and approximate differential privacy.

Lemma 8.2.13. *If P and Q are (ε, δ) -close, then for any $\beta > 0$, the sets*

$$A_+ := \{x : p(x) > e^{(1+\beta)\varepsilon} q(x)\} \quad \text{and} \quad A_- := \{x : p(x) \leq e^{-(1+\beta)\varepsilon} q(x)\}$$

satisfy

$$\max\{P(A_+), Q(A_-)\} \leq \frac{e^{\beta\varepsilon}\delta}{e^{\beta\varepsilon} - 1}, \quad \max\{P(A_-), Q(A_+)\} \leq \frac{e^{-\varepsilon}\delta}{e^{\beta\varepsilon} - 1}.$$

Conversely, if P and Q are (ε, δ) -pointwise close, then

$$P(A) \leq e^\varepsilon Q(A) + \delta \quad \text{and} \quad Q(A) \leq e^\varepsilon P(A) + \delta$$

for all sets A .

Proof Let $A = A_+ = \{x : p(x) > e^{(1+\beta)\varepsilon} q(x)\}$. Then

$$P(A) \leq e^\varepsilon Q(A) + \delta \leq e^{-\beta\varepsilon} P(A) + \delta,$$

so that $P(A) \leq \frac{\delta}{1 - e^{-\beta\varepsilon}}$. Similarly,

$$Q(A) \leq e^{-(1+\beta)\varepsilon} P(A) \leq e^{-\beta\varepsilon} Q(A) + e^{-(1+\beta)\varepsilon} \delta,$$

so that $Q(A) \leq e^{-(1+\beta)\varepsilon}\delta/(1-e^{-\beta\varepsilon}) = e^{-\varepsilon}\delta/(e^{\beta\varepsilon}-1)$. The set A_- satisfies the symmetric properties.

For the converse result, let $B = \{x : e^{-\varepsilon}q(x) \leq p(x) \leq e^{\varepsilon}q(x)\}$. Then for any set A we have

$$P(A) = P(A \cap B) + P(A \cap B^c) \leq e^{\varepsilon}Q(A \cap B) + \delta \leq e^{\varepsilon}Q(A) + \delta,$$

and the same inequalities yield $Q(A) \leq e^{\varepsilon}P(A) + \delta$. \square

That is, (ε, δ) -close distributions are $(2\varepsilon, \frac{e^{\varepsilon}+e^{-\varepsilon}}{e^{\varepsilon}-1}\delta)$ -pointwise close, and (ε, δ) -pointwise close distributions are (ε, δ) -close.

A minor extension of this lemma (taking $\beta = 1$ and applying the lemma twice) yields the following result.

Lemma 8.2.14. *Let P_0, P_1, P_2 be distributions on a space \mathcal{X} , each (ε, δ) -close. Then for any i, j, k , $j \neq k$, the set*

$$A_{jk} := \left\{ x \in \mathcal{X} : \log \frac{p_j(x)}{p_k(x)} > 3\varepsilon \right\} \quad \text{satisfies} \quad P_i(A_{jk}) \leq C\delta \max\{\varepsilon^{-1}, 1\}$$

for a numerical constant $C \leq 2$.

With Lemma 8.2.14 in hand, we can prove Proposition 8.2.11:

Proof of Proposition 8.2.11 Let $P_n^{(0)} \in \mathcal{P}_n$ denote the “true” sample. Consider the three channels Q_0, Q_1, Q_2 , which represent the induced distributions of $M(P_n^{(0)})$, $M(P_n)$, and $M(P'_n)$, respectively. Then by Lemma 8.2.14, with probability at least $1 - 2\delta \max\{\varepsilon^{-1}, 1\}$, $Z \sim Q_0$ belongs to the set $A = \{z \in \mathcal{Z} \mid e^{-3\varepsilon}q_1(z) \leq q_2(z) \leq e^{3\varepsilon}q_1(z)\}$. Calculating the odds ratios immediately gives the result. \square

Finally, we provide the proof of Proposition 8.2.12.

Proof of Proposition 8.2.12 Let $r = \alpha - 1$ for shorthand, and let $p = \frac{\alpha}{r} = \frac{\alpha}{\alpha-1} > 1$ and $p_* = \frac{p}{p-1} = \alpha$ be its conjugate. As in the proof of Proposition 8.2.11, let Q_0, Q_1 , and Q_2 represent the distributions of $Z = M(P_n^{(0)})$, $Z = M(P_n)$, and $Z = M(P'_n)$, respectively. We apply Hölder’s inequality: letting q_i be the density of Q_i with respect to some base measure $d\mu$ —which we know must exist by definition of Rényi differential privacy—we have

$$\begin{aligned} \mathbb{E} \left[\left(\frac{\pi(P_n | Z)}{\pi(P'_n | Z)} \right)^r \right] &= \int \left(\frac{q_1(z)\pi(P_n)}{q_2(z)\pi(P'_n)} \right)^r q_0(z) d\mu \\ &= \left(\frac{\pi(P_n)}{\pi(P'_n)} \right)^r \int \left(\frac{q_1(z)}{q_2(z)} \right)^r \frac{q_0(z)}{q_2(z)} q_2(z) d\mu \\ &\leq \left(\frac{\pi(P_n)}{\pi(P'_n)} \right)^r \left(\int \left(\frac{q_1(z)}{q_2(z)} \right)^{pr} q_2(z) d\mu \right)^{\frac{1}{p}} \left(\int \left(\frac{q_0(z)}{q_2(z)} \right)^{p_*} q_2(z) d\mu \right)^{\frac{1}{p_*}} \\ &= \left(\frac{\pi(P_n)}{\pi(P'_n)} \right)^r \exp \left(\frac{(\alpha-1)^2}{\alpha} D_\alpha(Q_1 \| Q_2) + \frac{\alpha-1}{\alpha} D_\alpha(Q_0 \| Q_2) \right) \\ &\leq \left(\frac{\pi(P_n)}{\pi(P'_n)} \right)^r \exp \left(\frac{(\alpha-1)^2 + \alpha - 1}{\alpha} \varepsilon \right) \end{aligned}$$

as $pr = \alpha$ and $p_* = \alpha$. Taking everything to the $1/(\alpha-1)$ power and gives the result. \square

8.3 Composition and privacy based on divergence

One of the major challenges in privacy is to understand what happens when a user participates in multiple studies, each providing different privacy guarantees. In this case, we might like to understand and control privacy losses even when the mechanisms for information release may depend on one another. Conveniently, all Rényi divergences provide strong guarantees on composition, essentially for free, and these then allow us to prove strong results on the composition of multiple private mechanisms.

8.3.1 Composition of Rényi-private channels

A natural idea to address composition is to attempt to generalize our chain rules for KL-divergence and related ideas to Rényi divergences. Unfortunately, this plan of attack does not quite work, as there is no generally accepted definition of a conditional Rényi divergence, and associated chain rules do not sum naturally. In situations in which individual divergence of associated elements of a joint distribution have bounded Rényi divergence, however, we can provide some natural bounds.

Indeed, consider the following essentially arbitrary scheme for data generation: we have distributions P and Q on a space \mathcal{Z}^n , where $Z_1^n \sim P$ and $Z_1^n \sim Q$ may exhibit arbitrary dependence. If, however, we can bound the conditional Rényi divergence between $P(Z_i | Z_1^{i-1})$ and $Q(Z_i | Z_1^{i-1})$, we can provide some natural tensorization guarantees. To set notation, let $P_i(\cdot | z_1^{i-1})$ be the (regular) conditional probability of Z_i conditional on $Z_1^{i-1} = z_1^{i-1}$ under P , and similarly for Q_i . We have the following theorem.

Theorem 8.3.1. *Let the conditions above hold, $\varepsilon_i < \infty$ for $i = 1, \dots, n$, and $\alpha \in [1, \infty]$. Assume that conditional on z_1^{i-1} , we have $D_\alpha(P_i(\cdot | z_1^{i-1}) \| Q_i(\cdot | z_1^{i-1})) \leq \varepsilon_i$. Then*

$$D_\alpha(P \| Q) \leq \sum_{i=1}^n \varepsilon_i.$$

Proof We assume without loss of generality that the conditional distributions $P_i(\cdot | z_1^{i-1})$ and Q_i are absolutely continuous with respect to a base measure μ on \mathcal{Z} .¹ Then we have

$$\begin{aligned} D_\alpha(P \| Q) &= \frac{1}{\alpha - 1} \log \int \prod_{i=1}^n \left(\frac{p_i(z_i | z_1^{i-1})}{q_i(z_i | z_1^{i-1})} \right)^\alpha q_i(z_i | z_1^{i-1}) d\mu^n(z_1^n) \\ &= \frac{1}{\alpha - 1} \log \int_{\mathcal{Z}_1^{n-1}} \left[\int \left(\frac{p_n(z_n | z_1^{n-1})}{q_n(z_n | z_1^{n-1})} \right)^\alpha q_n(z_n | z_1^{n-1}) d\mu(z_n) \right] \prod_{i=1}^{n-1} \left(\frac{p_i}{q_i} \right)^\alpha q_i d\mu^{n-1} \\ &\leq \frac{1}{\alpha - 1} \log \int_{\mathcal{Z}_1^{n-1}} \exp((\alpha - 1)\varepsilon_n) \prod_{i=1}^{n-1} \left(\frac{p_i(z_i | z_1^{i-1})}{q_i(z_i | z_1^{i-1})} \right)^\alpha q_i(z_i | z_1^{i-1}) d\mu^{n-1}(z_1^{n-1}) \\ &= \varepsilon_n + D_\alpha(P_1^{n-1} \| Q_1^{n-1}). \end{aligned}$$

Applying the obvious inductive argument then gives the result. \square

¹This is no loss of generality, as the general definition of f -divergences as suprema over finite partitions, or quantizations, of each X_i and Y_i separately, as in our discussion of KL-divergence in Chapter 2.2.2. Thus we may assume \mathcal{Z} is discrete and μ is a counting measure.

8.3.2 Privacy games and composition

To understand arbitrary composition of private channels, let us consider a privacy “game,” where an adversary may sequentially choose a dataset—in an arbitrary way—and then observes a private release Z_i of some mechanism applied to the dataset and the dataset with one entry (observation) modified. The adversary may then select a new dataset, and repeat the game. We then ask whether the resulting sequence of (private) observations Z_1^k remains private. Figure 8.1 captures this in an algorithmic form. Letting $Z_i^{(b)}$ denote the random observations under the bit $b \in \{0, 1\}$, whether

Input: Family of channels \mathcal{Q} and bit $b \in \{0, 1\}$.

Repeat: for $k = 1, 2, \dots$

- i. Adversary chooses arbitrary space \mathcal{X} , $n \in \mathbb{N}$, and two datasets $x^{(0)}, x^{(1)} \in \mathcal{X}^n$ with $d_{\text{ham}}(x^{(0)}, x^{(1)}) \leq 1$.
- ii. Adversary chooses private channel $Q_k \in \mathcal{Q}$.
- iii. Adversary observes one sample $Z_k \sim Q_k(\cdot \mid x^{(b)})$.

Figure 8.1. The privacy game. In this game, the adversary may *not* directly observe the private $b \in \{0, 1\}$.

the distributions of $(Z_1^{(0)}, \dots, Z_k^{(0)})$ and $(Z_1^{(1)}, \dots, Z_k^{(1)})$ are substantially different. Note that, in the game in Fig. 8.1, the adversary may track everything, and even chooses the mechanisms Q_k .

Now, let $Z^{(0)} = (Z_1^{(0)}, \dots, Z_k^{(0)})$ and $Z^{(1)} = (Z_1^{(1)}, \dots, Z_k^{(1)})$ be the outputs of the privacy game above, and let their respective marginal distributions be $Q^{(0)}$ and $Q^{(1)}$. We then make the following definition.

Definition 8.7. Let $\varepsilon \geq 0$, $\alpha \in [1, \infty]$, and $k \in \mathbb{N}$.

- (i) A collection \mathcal{Q} of channels satisfies (ε, α) -Rényi privacy under k -fold adaptive composition if, in the privacy game in Figure 8.1, the distributions $Q^{(0)}$ and $Q^{(1)}$ on $Z^{(0)}$ and $Z^{(1)}$, respectively, satisfy $D_\alpha(Q^{(0)} \| Q^{(1)}) \leq \varepsilon$ and $D_\alpha(Q^{(1)} \| Q^{(0)}) \leq \varepsilon$.
- (ii) Let $\delta > 0$. Then a collection \mathcal{Q} of channels satisfies (ε, δ) -differential privacy under k -fold adaptive composition if $D_\infty^\delta(Q^{(0)} \| Q^{(1)}) \leq \varepsilon$ and $D_\infty^\delta(Q^{(1)} \| Q^{(0)}) \leq \varepsilon$.

By considering a special case centered around a particular individual in the game 8.1, we can gain some intuition for the definition. Indeed, suppose that an individual has some data x_0 ; in each round of the game the adversary generates two datasets, one containing x_0 and the other identical except that x_0 is removed. Then satisfying Definition 8.7 captures the intuition that an individual’s privacy remains protected, even in the face of multiple (private) accesses of the individual’s data.

As an immediate corollary to Theorem 8.3.1, we then have the following.

Corollary 8.3.2. Assume that each channel in the game in Fig. 8.1 is (ε_i, α) -Rényi private. Then the arbitrary composition of k such channels remains $(\sum_{i=1}^k \varepsilon_i, \alpha)$ -Rényi private.

More sophisticated corollaries are possible once we start to use the connections between privacy measures we outline in Section 8.2.2. In this case, we can develop so-called *advanced composition* rules, which sometimes suggest that privacy degrades more slowly than might be expected under adaptive composition.

Corollary 8.3.3. *Assume that each channel in the game in Fig. 8.1 is ε -differentially private. Then the composition of k such channels is $k\varepsilon$ -differentially private. Additionally, the composition of k such channels is*

$$\left(\frac{3k}{2}\varepsilon^2 + \sqrt{6k \log \frac{1}{\delta}} \cdot \varepsilon, \delta \right)$$

differentially private for all $\delta > 0$.

Proof The first claim is immediate: for $Q^{(0)}, Q^{(1)}$ as in Definition 8.7, we know that $D_\alpha(Q^{(0)} \| Q^{(1)}) \leq k\varepsilon$ for all $\alpha \in [1, \infty]$ by Theorem 8.3.1 coupled with Proposition 8.2.5 (or Corollary 8.2.6).

For the second claim, we require a bit more work. Here, we use the bound $\frac{3\alpha}{2}\varepsilon^2$ in the Rényi privacy bound in Corollary 8.2.6. Then we have for any $\alpha \geq 1$ that

$$D_\alpha(Q^{(0)} \| Q^{(1)}) \leq \frac{3k\alpha}{2}\varepsilon^2$$

by Theorem 8.3.1. Now we apply Proposition 8.2.7 and Corollary 8.2.8, which allow us to conclude (ε, δ) -differential privacy from Rényi privacy. Indeed, by the preceding display, setting $\eta = 1 + \alpha$, we have that the composition is $(\frac{3k}{2}\varepsilon^2 + \frac{3k\eta}{2}\varepsilon^2 + \frac{1}{\eta} \log \frac{1}{\delta}, \delta)$ -differentially private for all $\eta > 0$ and $\delta > 0$. Optimizing over η gives the second result. \square

We note in passing that it is possible to get slightly sharper results than those in Corollary 8.3.3; indeed, using ideas from Exercise 4.4 it is possible to achieve $(k\varepsilon(e^\varepsilon - 1) + \sqrt{2k \log \frac{1}{\delta}}\varepsilon, \delta)$ -differential privacy under adaptive composition.

A more sophisticated result, which shows adaptive composition for (ε, δ) -differentially private channels, is also possible using Lemma 8.2.10.

Theorem 8.3.4. *Assume that each channel in the game in Fig. 8.1 is (ε, δ) -differentially private. Then the composition of k such channels is $(k\varepsilon, k\delta)$ -differentially private. Additionally, they are*

$$\left(\frac{3k}{2}\varepsilon^2 + \sqrt{6k \log \frac{1}{\delta_0}} \cdot \varepsilon, \delta_0 + \frac{k\delta}{1 + e^\varepsilon} \right)$$

differentially private for all $\delta_0 > 0$.

Proof Consider the channels Q_i in Fig. 8.1. As each satisfies $D_\infty(Q_i(\cdot | x^{(0)}) \| Q_i(\cdot | x^{(1)})) \leq \varepsilon$ and $D_\infty(Q_i(\cdot | x^{(1)}) \| Q_i(\cdot | x^{(0)})) \leq \varepsilon$, Lemma 8.2.10 guarantees the existence (at each sequential step, which may depend on the preceding $i - 1$ outputs) of probability measures $Q_i^{(0)}$ and $Q_i^{(1)}$ such that $D_\infty(Q_i^{(1-b)} \| Q_i^{(b)}) \leq \varepsilon$, $\|Q_i^{(b)} - Q_i(\cdot | x^{(b)})\|_{\text{TV}} \leq \delta/(1 + e^\varepsilon)$ for $b \in \{0, 1\}$.

Note that by construction (and Theorem 8.3.1) we have $D_\alpha(Q_1^{(b)} \cdots Q_k^{(b)} \| Q_1^{(1-b)} \cdots Q_k^{(1-b)}) \leq \min\{\frac{3k\alpha}{2}\varepsilon^2, k\varepsilon\}$, where $Q^{(b)}$ denotes the joint distribution on Z_1, \dots, Z_k under bit b . We also have by the triangle inequality that $\|Q_1^{(b)} \cdots Q_k^{(b)} - Q^{(b)}\|_{\text{TV}} \leq k\delta/(1 + e^\varepsilon)$ for $b \in \{0, 1\}$. (See Exercise 2.16.) As a consequence, we see as in the proof of Corollary 8.3.3 that the composition is $(\frac{3k}{2}\varepsilon^2 + \frac{3k\eta}{2}\varepsilon^2 + \frac{1}{\eta} \log \frac{1}{\delta_0}, \delta_0 + k\delta/(1 + e^\varepsilon))$ -differentially private for all $\eta > 0$ and δ_0 . Optimizing gives the result. \square

As a consequence of these results, we see that whenever the privacy parameter $\varepsilon < 1$, it is possible to compose multiple privacy mechanisms together and have privacy penalty scaling only as the worse of $\sqrt{k\varepsilon}$ and $k\varepsilon^2$, which is substantially better than the “naive” bound of $k\varepsilon$. Of course, a challenge here—relatively unfrequently discussed in the privacy literature—is that when $\varepsilon \geq 1$, which is a frequent case for practical deployments of privacy, all of these bounds are much worse than a naive bound that k -fold composition of ε -differentially private algorithms is $k\varepsilon$ -differentially private.

8.4 Additional mechanisms and privacy-preserving algorithms

Since the introduction of differential privacy, a substantial literature has grown providing mechanisms for different estimation, learning, and data release problems. Here, we describe a few of those beyond the basic noise addition schemes we have thus far developed, highlighting a few applications along the way. One major challenge with the naive approaches is that they rely on *global* sensitivity of the functions to be estimated, rather than local sensitivities—a worst case notion that sometimes forces privacy to add unnecessary noise. In Section 8.4.2, we give one potential approach to this problem, which we develop further in exercises and revisit in optimality guarantees in sequential chapters. Our view is necessarily somewhat narrow, but the results here can form a natural starting point for further work in this area.

8.4.1 The exponential mechanism

In many statistical, learning, and other problems, there is a natural notion of loss (or conversely, utility) in releasing a potentially noisy result of some computation. We abstract this by considering the input space \mathcal{P}_n of samples of size n (that is, empirical distributions) and output space \mathcal{Z} along with a loss function $\ell : \mathcal{P}_n \times \mathcal{Z} \rightarrow \mathbb{R}$, where $\ell(P_n, z)$ measures the loss of z on an input $P_n \in \mathcal{P}_n$. For example, if we wish to compute a function $f : \mathcal{P}_n \rightarrow \mathbb{R}$, a natural notion of loss is $\ell(P_n, z) = |f(P_n) - z|$ for $z \in \mathbb{R}$. As a more sophisticated and somewhat abstract formulation, suppose we wish to release a sample distribution \tilde{P} approximating an input sample $P_n \in \mathcal{P}_n$, where we wish \tilde{P} to be accurate for most statistical queries in some family, that is, $\frac{1}{n} \sum_{i=1}^n \phi(x_i) \approx \mathbb{E}_{\tilde{P}}[\phi(X)]$ for all $\phi \in \Phi$. Then a natural loss is $\ell(P_n, \tilde{P}) = \sup_{\phi \in \Phi} |\mathbb{E}_{P_n} \phi(X) - \mathbb{E}_{\tilde{P}}[\phi(X)]|$.

In scenarios in which we have such a loss, the abstract *exponential mechanism* provides an attractive approach. We assume that for each $z \in \mathcal{Z}$, the loss $\ell(\cdot, z)$ has (global) sensitivity L , i.e., $|\ell(P_n, z) - \ell(P'_n, z)| \leq L$ for all neighboring $P_n, P'_n \in \mathcal{P}_n$. We assume we have a base measure μ on \mathcal{Z} , and then define the exponential mechanism by

$$\mathbb{P}(M(P_n) \in A) = \frac{1}{\int \exp(-\frac{\varepsilon}{L}\ell(P_n, z))d\mu(z)} \int_A \exp\left(-\frac{\varepsilon}{L}\ell(P_n, z)\right) d\mu(z), \quad (8.4.1)$$

assuming $\int e^{-\frac{\varepsilon}{L}\ell(x, z)}d\mu(z)$ is finite for each $P_n \in \mathcal{P}_n$. (Typically, one assumes ℓ takes on values in \mathbb{R}_+ and μ is a finite measure, making the last assumption trivial.) That is, the exponential mechanism M releases $Z = M(P_n)$ with probability proportional to

$$\exp\left(-\frac{\varepsilon}{L}\ell(P_n, z)\right).$$

That the mechanism (8.4.1) is 2ε -differentially private is immediate: for any neighboring P_n, P'_n ,

we have

$$\begin{aligned} \frac{Q(A \mid P_n)}{Q(A \mid P'_n)} &= \frac{\int \exp(-\frac{\varepsilon}{L} \ell(P'_n, z)) d\mu(z)}{\int \exp(-\frac{\varepsilon}{L} \ell(P_n, z)) d\mu(z)} \cdot \frac{\int_A \exp(-\frac{\varepsilon}{L} \ell(P_n, z)) d\mu(z)}{\int_A \exp(-\frac{\varepsilon}{L} \ell(P'_n, z)) d\mu(z)} \\ &\leq \sup_{z \in \mathcal{Z}} \left\{ \exp \left(\frac{\varepsilon}{L} [\ell(P_n, z) - \ell(P'_n, z)] \right) \right\} \cdot \sup_{z \in A} \left\{ \exp \left(\frac{\varepsilon}{L} [\ell(P'_n, z) - \ell(P_n, z)] \right) \right\} \leq \exp(2\varepsilon). \end{aligned}$$

As a first (somewhat trivial) example, we can recover the Laplace mechanism:

Example 8.4.1 (The Laplace mechanism): We can recover Example 8.1.3 through the exponential mechanism. Indeed, suppose that we wish to release $f : \mathcal{P}_n \rightarrow \mathbb{R}^d$, where $\text{GS}_1(f) \leq L$. Then taking $z \in \mathbb{R}^d$, $\ell(P_n, z) = \|f(P_n) - z\|_1$, and μ to be the usual Lebesgue measure on \mathbb{R}^d , the exponential mechanism simply uses density

$$q(z \mid P_n) \propto \exp \left(-\frac{\varepsilon}{L} \|f(P_n) - z\|_1 \right),$$

which is the Laplace mechanism. \diamond

One challenge with the exponential mechanism (8.4.1) is that it is somewhat abstract and is often hard to compute, as it requires evaluating an often high-dimensional integral to sample from. Yet it provides a nice abstract mechanism with strong privacy guarantees and, as we shall see, good utility guarantees. For the moment, we defer further examples and provide utility guarantees when $\mu(\mathcal{Z})$ is finite, giving bounds based on the measure of “bad” solutions. For notational convenience, we define the optimal value

$$\ell^*(P_n) = \inf_{z \in \mathcal{Z}} \ell(P_n, z),$$

assuming tacitly that it is finite, and the sublevel sets

$$S_t := \{z \in \mathcal{Z} \mid \ell(P_n, z) \leq \ell^*(P_n) + t\}.$$

With these definitions, we have the following proposition.

Proposition 8.4.2. *Let $t \geq 0$. Then for the exponential mechanism (8.4.1), if $Z \sim Q(\cdot \mid P_n)$ then*

$$\ell(P_n, Z) \leq \ell^*(P_n) + 2t$$

with probability at least $1 - \exp \left(-\frac{\varepsilon t}{L} + \log \frac{\mu(\mathcal{Z})}{\mu(S_t)} \right)$.

Proof Assume without loss of generality (by scaling) that the global Lipschitzian (sensitivity) constant of ℓ is $L = 1$. Then for $Z \sim Q(\cdot \mid P_n)$, we have

$$\begin{aligned} \mathbb{P}(\ell(P_n, Z) \geq \ell^*(P_n) + 2t) &= \frac{\int_{S_{2t}^c} \exp(-\varepsilon \ell(P_n, z)) d\mu(z)}{\int \exp(-\varepsilon \ell(P_n, z)) d\mu(z)} = \frac{\int_{S_{2t}^c} \exp(-\varepsilon (\ell(P_n, z) - \ell^*(P_n))) d\mu(z)}{\int \exp(-\varepsilon (\ell(P_n, z) - \ell^*(P_n))) d\mu(z)} \\ &\leq \frac{\int_{S_{2t}^c} \exp(-2\varepsilon t) d\mu(z)}{\int_{S_t} \exp(-\varepsilon (\ell(P_n, z) - \ell^*(P_n))) d\mu(z)} \leq \exp(-\varepsilon t) \frac{\mu(S_{2t}^c)}{\mu(S_t)}, \end{aligned}$$

where the last inequality uses that $\ell(P_n, z) - \ell^*(P_n) \leq t$ on S_t . \square

We can provide a few simplifications of this result in different special cases. For example, if \mathcal{Z} is finite with cardinality $\text{card}(\mathcal{Z})$, then Proposition 8.4.2 implies that taking μ to be the counting measure on \mathcal{Z} we have

Corollary 8.4.3. *In addition to the conditions in Proposition 8.4.2, assume that $\text{card}(\mathcal{Z})$ is finite. Then for any $u \in (0, 1)$, with probability at least $1 - u$,*

$$\ell(P_n, Z) \leq \ell^*(P_n) + \frac{2L}{\varepsilon} \log \frac{\text{card}(\mathcal{Z})}{u}.$$

That is, with extremely high probability, the loss of Z from the exponential mechanism is at most logarithmic in $\text{card}(\mathcal{Z})$ and grows only linearly with the global sensitivity L .

A second corollary allows us to bound the expected loss of the exponential mechanism, assuming we have some control over the measure of the sublevel sets S_t .

Corollary 8.4.4. *Let $t \geq 0$ be the smallest scalar such that $t \geq \frac{2L}{\varepsilon} \log \frac{\mu(\mathcal{Z})}{\mu(S_t)}$ and $t \geq \frac{L}{\varepsilon}$. Then Z drawn from the exponential mechanism (8.4.1) satisfies*

$$\mathbb{E}[\ell(P_n, Z)] \leq \ell^*(P_n) + t + \frac{2L}{\varepsilon} \leq \ell^*(P_n) + 3t \leq \ell^*(P_n) + O(1) \frac{L}{\varepsilon} \log \left(1 + \frac{\mu(\mathcal{Z})}{\mu(S_t)} \right).$$

Proof We first recall that if $W \geq 0$ is a nonnegative random variable, then by a change of variables, $\mathbb{E}[W] = \int_0^\infty \mathbb{P}(W \geq t) dt$. Take $\ell(P_n, Z) - \ell^*(P_n) \geq 0$ as our random variable, fix any $t_0 \geq 0$, and let $\rho = \log \frac{\mu(\mathcal{Z})}{\mu(S_{t_0})}$. Then by Proposition 8.4.2 we have

$$\begin{aligned} \mathbb{E}[\ell(P_n, Z) - \ell^*(P_n)] &\leq t_0 + \int_{t_0}^\infty \mathbb{P}(\ell(P_n, Z) - \ell^*(P_n) \geq t) dt \\ &= t_0 + 2 \int_{t_0/2}^\infty \mathbb{P}(\ell(P_n, Z) - \ell^*(P_n) \geq 2t) dt \\ &\leq t_0 + 2 \int_{t_0/2}^\infty \exp \left(-\frac{\varepsilon t}{L} + \log \frac{\mu(\mathcal{Z})}{\mu(S_t)} \right) dt \\ &\leq t_0 + 2e^\rho \int_{t_0/2}^\infty \exp \left(-\frac{\varepsilon t}{L} \right) dt = t_0 + \frac{2L}{\varepsilon} \exp \left(\rho - \frac{\varepsilon t_0}{2L} \right). \end{aligned}$$

Take t_0 as in the statement of the corollary to obtain the result. \square

Corollary 8.4.4 may seem a bit circular: we require the ratio $\mu(\mathcal{Z})/\mu(S_t)$ to be controlled—but it is relatively straightforward to use it (and Proposition 8.4.2) with a bit of care and standard bounds on volumes.

Example 8.4.5 (Empirical risk minimization via the exponential mechanism): We consider the empirical risk minimization problem, where we have losses $\ell : \Theta \times \mathcal{X} \rightarrow \mathbb{R}_+$, where $\Theta \subset \mathbb{R}^d$ is a parameter space of interest, and we wish to choose

$$\hat{\theta}_n \in \operatorname{argmin}_{\theta \in \Theta} \left\{ L(\theta, P_n) := \frac{1}{n} \sum_{i=1}^n \ell(\theta, x_i) \right\}$$

where $P_n = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i}$. We make a few standard assumptions: first, for simplicity, that n is large enough that $\frac{n}{d} \geq \varepsilon$. We also assume that $\Theta \subset \mathbb{R}^d$ is an ℓ_2 -ball of radius R , that $\theta \mapsto \ell(\theta, x_i)$ is M -Lipschitz for all x_i , and that $\ell(\theta, x_i) \in [0, 2MR]$ for all $\theta \in \Theta$. (Note that this last is no loss of generality, as $\ell(\theta, x_i) - \inf_{\theta \in \Theta} \ell(\theta, x_i) \leq M \sup_{\theta, \theta' \in \Theta} \|\theta - \theta'\|_2 \leq 2MR$.)

Take the empirical loss $L(\theta, P_n)$ as our criterion function for the exponential mechanism, which evidently satisfies $|L(\theta, P_n) - L(\theta, P'_n)| \leq \frac{2MR}{n}$ whenever $d_{\text{ham}}(P_n, P'_n) \leq 1$, so that we release θ with density

$$q(\theta | x) \propto \exp\left(-\frac{n\varepsilon}{2MR}L(\theta, P_n)\right).$$

Let $\hat{\theta}_n$ be the empirical minimizer as above; then by the Lipschitz continuity of ℓ , the sublevel set S_t evidently satisfies

$$S_t \supset \left\{ \theta \in \Theta \mid \|\theta - \hat{\theta}_n\|_2 \leq \frac{t}{M} \right\}.$$

Then a volume calculation (with the factor of 2 necessary because we may have $\hat{\theta}_n$ on the boundary of Θ) yields that for μ the Lebesgue measure,

$$\frac{\mu(S_t)}{\mu(\mathcal{Z})} \geq \left(\frac{t}{2MR}\right)^d.$$

As a consequence, by Corollary 8.4.4, whenever $t \geq O(1)\frac{MR}{n\varepsilon} \cdot d \log \frac{MR}{t}$, we have $\mathbb{E}[L(\theta, P_n) | P_n] \leq L(\hat{\theta}_n, P_n) + 3t$. The choice $t = O(1)\frac{MRd}{n\varepsilon}$ suffices whenever $\frac{\varepsilon}{d} \leq 1$, so we obtain

$$\mathbb{E}[L(\theta, P_n)] \leq L(\hat{\theta}_n, P_n) + O(1)\frac{MRd}{n\varepsilon} \log \frac{n\varepsilon}{d},$$

whenever $\frac{d}{n\varepsilon} \leq 1$. Notably, standard empirical risk minimization (recall Chapter 5.2) typically achieves rates of convergence roughly of MR/\sqrt{n} , so that the gap of the exponential mechanism is lower order whenever $\frac{d}{\sqrt{n\varepsilon}} \leq 1$. \diamond

8.4.2 Local sensitivities and the inverse sensitivity mechanism

A particular choice of the exponential mechanism (8.4.1) can provide strong optimality guarantees for 1-dimensional quantities, and appears to be the “right” mechanism (in principle) when one wishes to estimate a scalar-valued functional $f(P_n)$. A better (in principle) algorithm than noise addition schemes using the global sensitivity $\text{GS}(f) = \sup |f(P_n) - f(P'_n)|$ is to use a *local* notion of sensitivity: we are only concerned with adding noise commensurate with the changes of f near $P_n \in \mathcal{P}_n$. With this in mind, define the *modulus of continuity* of f at P_n by

$$\omega_f(k; P_n) := \sup \{|f(P'_n) - f(P_n)| \mid d_{\text{ham}}(P_n, P'_n) \leq k\},$$

which measures the amount that changing k observations in P_n can change the function f . In the privacy literature, the particular choice $k = 1$ yields the *local sensitivity*

$$\text{LS}(f, P_n) := \sup \{|f(P'_n) - f(P_n)| \mid d_{\text{ham}}(P'_n, P_n) = 1\} = \omega_f(1; P_n). \quad (8.4.2)$$

A naive strategy, then, would be to release

$$Z = f(P_n) + \frac{\text{LS}(f, P_n)}{\varepsilon} \cdot W \quad \text{for } W \sim \text{Laplace}(1),$$

which is analogous to the Laplace mechanism (8.1.3), except that the noise scales with the local sensitivity of f at P_n . The issue, as the next example makes clear, is that the scale of this noise can compromise privacy.

Example 8.4.6 (The sensitivity of the sensitivity): Consider estimating a median $f(P_n) = \text{med}(P_n)$, where the data $x \in [0, 1]$, where $n = 2m + 1$ for simplicity, to make the median unique. If the sample consists of m points $x_i = 0$ and $m + 1$ points $x_i = 1$, then the sensitivity $\omega_f(1, P_n) = 1$, the maximal value—we simply move one example from $x_i = 1$ to $x_i = 0$, changing the median from $\text{med}(P_n) = 1$ to 0. On the other hand, on the sample P'_n with $m - 1$ points $x_i = 0$ and $m + 2$ points $x_i = 1$, the sensitivity $\omega_f(1, P'_n) = 0$, because changing a single example cannot move the median from $f(P'_n) = 1$. \diamond

Instead of using the inherently unstable quantity ω , then, we can instead use, essentially, its inverse: define the *inverse sensitivity*

$$d_f(t, P_n) := \inf \{d_{\text{ham}}(P'_n, P_n) \mid f(P'_n) = t\}, \quad (8.4.3)$$

where $d_f(t, P_n) = +\infty$ if no P'_n yields $f(P'_n) = t$. So $d_f(t, P_n)$ counts the number of examples that must be changed in the sample P_n to move $f(P_n)$ to a target t , and by inspection, always satisfies

$$|d_f(t, P_n) - d_f(t, P'_n)| \leq d_{\text{ham}}(P_n, P'_n).$$

Then the *inverse sensitivity mechanism* releases a value t with probability density proportional to

$$q(t \mid P_n) \propto \exp\left(-\frac{\varepsilon}{2} d_f(t, P_n)\right). \quad (8.4.4)$$

Implicit in the definition (8.4.4) is a base measure μ , typically one of Lebesgue measure or counting measure on a discrete set. Then a quick calculation (or recognition that the density (8.4.4) is a particular instance of the exponential mechanism) gives the following proposition.

Proposition 8.4.7. *Let M be the inverse sensitivity mechanism with density (8.4.4). Then M is ε -differentially private.*

As in the general exponential mechanism (8.4.1), efficiently sampling from the density (8.4.4) can be challenging. Some cases admit easier reformulations.

Example 8.4.8 (Mean estimation with bounded data): Suppose the data $x \in [a, b]$ are bounded and we wish to estimate the sample mean $f(P_n) = \mathbb{E}_{P_n}[X] = \bar{x}_n$, where $P_n = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i}$. Changing a single observation can move the mean by at most $\frac{b-a}{n}$ (replace $x_i = a$ with $x'_i = b$). Thus, while discretization issues and that we may have $x_i \notin \{a, b\}$ make precisely computing d_f tedious, the approximation

$$d_{\text{mean}}(t, P_n) = \left\lceil \frac{n|t - \bar{x}_n|}{b - a} \right\rceil,$$

where we define $d_{\text{mean}}(t, P_n) = +\infty$ for $t \notin [a, b]$, is both Lipschitz (with respect to the Hamming metric) in the sample P_n , and approximates $d_f(t, P_n)$. (See Exercise 8.8 for a more general approach justifying this particular approximation.) The approximation

$$q(t \mid P_n) := \frac{\exp(-\frac{\varepsilon}{2} d_{\text{mean}}(t, P_n))}{\int_a^b \exp(-\frac{\varepsilon}{2} d_{\text{mean}}(s, P_n)) ds} \quad (8.4.5)$$

to the density (8.4.4) is thus ε -differentially private,

The density (8.4.5) yields a particular step-like density. Define the shells

$$S_k = \left\{ \left[\bar{x}_n - k \frac{b-a}{n}, \bar{x}_n - (k-1) \frac{b-a}{n} \right] \cup \left[\bar{x}_n + (k-1) \frac{b-a}{n}, \bar{x}_n + k \frac{b-a}{n} \right] \right\} \cap [a, b]$$

corresponding to the amount the mean may change if we modify k examples and let $\text{Vol}(S_k)$ be volume (length) of the intervals making up S_k . To sample from the density (8.4.5), note that the denominator $C(P_n) := \int_a^b \exp(-\frac{\varepsilon}{2} d_{\text{mean}}(s, P_n)) ds = \sum_{k=1}^n \text{Vol}(S_k) e^{-\frac{k\varepsilon}{2}}$. Then we draw an index $I \in [n]$ with probability $\mathbb{P}(I = k) = \text{Vol}(S_k) e^{-\varepsilon k/2} / C(P_n)$, and then choose t uniformly at random within S_k . \diamond

Example 8.4.9 (Median estimation): For the median, the inverse sensitivity takes a particularly clean form, making sampling from the density (8.4.4) fairly straightforward. In this case, for a sample $P_n = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i}$, where $x_i \in \mathbb{R}$, we have

$$d_f(t, P_n) = \text{card} \{i \in [n] \mid x_i \in [f(P_n), t]\},$$

the number of examples between the median $f(P_n)$ and putative target t . If the data lie in a range $x \in [a, b]$, then the density q is relatively straightforward to compute. Similar to the approach to the stepped density in Example 8.4.8, divide $[a, b]$ into the intervals

$$S_k^- := [a_k^-, a_{k-1}^-] \quad \text{and} \quad S_k^+ := [a_{k-1}^+, a_k^+], \quad k = 1, \dots, n/2,$$

where

$$a_k^- = \inf \{f(P'_n) \mid d_{\text{ham}}(P'_n, P_n) \leq k\} \quad \text{and} \quad a_k^+ = \sup \{f(P'_n) \mid d_{\text{ham}}(P'_n, P_n) \leq k\}.$$

That is, a_k^- is the smallest we can make the median by changing k examples and a_k^+ the largest, corresponding to the $\frac{1}{2} - \frac{k}{n}$ and $\frac{1}{2} + \frac{k}{n}$ quantiles of the sample P_n , where the 0 quantile is a and 1 quantile is b . Then defining the normalization constant

$$C(P_n) := \int_a^b \exp\left(-\frac{\varepsilon}{2} d_f(t, P_n)\right) dt = \sum_{k=1}^n \text{Vol}(S_k^- \cup S_k^+) \exp\left(-\frac{\varepsilon}{2} k\right)$$

(where the volume is simply interval length), we may sample from the density (8.4.4) by first drawing a random index $I \in \{1, \dots, n\}$ with probability proportional to

$$\mathbb{P}(I = k \mid P_n) = \frac{\text{Vol}(S_k^- \cup S_k^+)}{C(P_n)} \exp\left(-\frac{\varepsilon}{2} k\right),$$

then drawing t uniformly at random in the each of the intervals S_k^- or S_k^+ with probabilities $\text{Vol}(S_k^-) / \text{Vol}(S_k^- \cup S_k^+)$ or $\text{Vol}(S_k^+) / \text{Vol}(S_k^- \cup S_k^+)$, respectively. \diamond

The particular sampling strategies—where we construct concentric shells S_k around $f(P_n)$ and sample from these with geometrically decaying probabilities $e^{-k\varepsilon/2}$ —point toward more general sampling strategies and optimality guarantees for the inverse sensitivity mechanism. Define the “shells”

$$S_k := \{f(P'_n) \mid d_{\text{ham}}(P_n, P'_n) = k\}.$$

We focus on sampling from the density (8.4.4) in the case $t \in \mathbb{R}$, so sampling is equivalent to drawing an index $I \in [n]$ with probability

$$\mathbb{P}(I = k \mid P_n) = \frac{1}{C(P_n)} e^{-\frac{\varepsilon}{2}k} \quad \text{for} \quad C(P_n) := \sum_{k=1}^n \text{Vol}(S_k) e^{-\frac{\varepsilon}{2}k}, \quad (8.4.6)$$

then choosing t uniformly at random in S_k .

Define the shorthand $\omega(k) = \omega_f(k, P_n)$. Then the values $t \in S_k$ all satisfy $|f(P_n) - t| \leq \omega(k)$, and so the inverse sensitivity mechanism M guarantees

$$\mathbb{E}[|M(P_n) - f(P_n)|] \leq \sum_{k=1}^n \mathbb{P}(M(P_n) \in S_k) \omega(k).$$

Now our calculations become heuristic, where we make an effort to give the rough flavor of results possible, and later apply the care necessary for tighter guarantees. Suppose that the interval lengths $\text{Vol}(S_k)$ are of the same order for $k \lesssim \frac{1}{\varepsilon}$, and grow only polynomially quickly for $k \gg \frac{1}{\varepsilon}$. Then we have the heuristic bound $C(P_n) := \sum_{k=1}^n \text{Vol}(S_k) e^{-k\varepsilon/2} \gtrsim \text{Vol}(S_1) \sum_{k=1}^n e^{-k\varepsilon/2} \gtrsim \varepsilon^{-1} \text{Vol}(S_1)$, while

$$\mathbb{E}[|M(P_n) - f(P_n)|] \leq \sum_{k=1}^n \frac{\text{Vol}(S_k) e^{-k\varepsilon/2}}{\sum_{i=1}^n \text{Vol}(S_i) e^{-i\varepsilon/2}} \omega(k) \stackrel{\text{heuristic}}{\lesssim} \sum_{k=1}^n \varepsilon e^{-k\varepsilon/2} \omega(k) \lesssim \max_k e^{-k\varepsilon/2} \omega(k),$$

where the heuristic inequality is our bound on the normalizing constant $C(P_n)$, and the final bound follows because maxima are larger than (weighted) averages. Continuing the heuristic derivation, the final maximum places exponentially small weight on $\omega(k)$ for $k \gg \frac{1}{\varepsilon}$. Thus—and again, this is non-rigorous—we expect roughly that

$$\mathbb{E}[|M(P_n) - f(P_n)|] \stackrel{\text{heuristic}}{\lesssim} \max_k e^{-k\varepsilon/2} \omega(k) \stackrel{\text{heuristic}}{\lesssim} \omega_f\left(\frac{c}{\varepsilon}, P_n\right), \quad (8.4.7)$$

where c is some numerical constant.

To gain some intuition for the claims of optimality we have made, let us revisit the equivalent definitions of privacy that repose on testing, as in Eq. (8.1.4) and Proposition 8.1.6. By the definition of differential privacy, the inverse sensitivity mechanism satisfies

$$\mathbb{P}(M(P_n) \in A) \leq e^{k\varepsilon} \mathbb{P}(M(P'_n) \in A)$$

for any samples P_n, P'_n satisfying $d_{\text{ham}}(P_n, P'_n) \leq k$. So for $k \leq \frac{1}{\varepsilon}$, we have

$$\mathbb{P}(M(P_n) \in A) \leq \exp(1) \mathbb{P}(M(P'_n) \in A),$$

and so no procedure exists that can test whether the sample is P_n or P'_n with probability of error less than e^{-2} , by Proposition 8.1.6. Thus, at a fundamental level, *no* procedure can reliably distinguish the outputs of $M(P_n)$ from those of $M(P'_n)$ when P_n and P'_n differ in only $1/\varepsilon$ examples. Thus, we cannot expect to estimate $f(P_n)$ to accuracy better than $\omega_f(\frac{1}{\varepsilon}, P_n)$, and so for any ε -differentially private mechanism M and P_n , there exists $P'_n \in \mathcal{P}_n$ with $d_{\text{ham}}(P_n, P'_n) \leq \frac{1}{\varepsilon}$ and for which

$$\max_{\hat{P} \in \{P_n, P'_n\}} \mathbb{E}[|M(\hat{P}) - f(\hat{P})|] \gtrsim \omega_f\left(\frac{1}{\varepsilon}, P_n\right), \quad (8.4.8)$$

which the heuristic calculation (8.4.7) achieves.

To provide more rigorous guarantees requires restrictions on the functions f whose values we wish to release. The simplest is that the function $f : \mathcal{P}_n \rightarrow \mathbb{R}$ obey a natural ordering property, where larger changes in the sample distribution P_n beget larger changes in f .

Definition 8.8. A function $f : \mathcal{P}_n \rightarrow \mathbb{R}$ is sample monotone if for each $s, t \in f(\mathcal{P}_n)$ satisfying $f(P_n) \leq s \leq t$ or $t \leq s \leq f(P_n)$, we have $d_f(s, P_n) \leq d_f(t, P_n)$.

The mean and median (Examples 8.4.8 and 8.4.9) are both sample monotone. So, too, are appropriately continuous functions f . For this, we make the obvious identification of $f : \mathcal{P}_n \rightarrow \mathbb{R}$ with the induced function on \mathcal{X}^n by defining $f_{\mathcal{X}}(x_1^n) := f(n^{-1} \sum_{i=1}^n \mathbf{1}_{x_i})$. Then we say $f : \mathcal{P}_n \rightarrow \mathbb{R}$ is continuous if the induced function $f_{\mathcal{X}}$ is.

Observation 8.4.10. Let $f : \mathcal{P}_n \rightarrow \mathbb{R}$ be continuous and \mathcal{X} convex. Then f is sample monotone.

Proof Identify f with its induced function $f_{\mathcal{X}}$ for notational simplicity, and let $x \in \mathcal{X}^n$, $f(x) \leq s \leq t$, and $P_n = n^{-1} \sum_{i=1}^n \mathbf{1}_{x_i}$ be the empirical distribution associated with x . We show that $d_f(s, P_n) \leq d_f(t, P_n)$. If $d_f(t, P_n) = +\infty$, then the desired inequality holds trivially. Otherwise, let $x' \in \mathcal{X}^n$ satisfy $f(x') = t$ and $d_{\text{ham}}(x, x') = d_f(t, P_n)$. Then the function $g(\lambda) := f((1 - \lambda)x + \lambda x')$ is continuous in λ and satisfies $g(0) = f(x) \leq g(1) = f(x') = t$. By the intermediate value theorem, there exists $\lambda_s \in [0, 1]$ with $g(\lambda_s) = s$, and as \mathcal{X} is convex the vector $x_s = (1 - \lambda_s)x + \lambda_s x' \in \mathcal{X}^n$ satisfies $f(x_s) = g(\lambda_s) = s$. That x_s is a convex combination of x and x' then implies $d_f(s, P_n) \leq d_{\text{ham}}(x, x_s) \leq d_{\text{ham}}(x, x') = d_f(t, P_n)$. \square

With Definition 8.8 in place, we can provide a few stronger guarantees for the inverse sensitivity mechanism. To avoid pathological sampling issues, one replaces the inverse sensitivity $d_f(t, P_n)$ with a “smoothed” version, where for $\rho \geq 0$ we define

$$d_{f,\rho}(t, P_n) := \inf \{d_{\text{ham}}(P_n, P'_n) \mid |f(P'_n) - t| \leq \rho\}.$$

(Pathological cases include estimating the median where the sample P_n consists of a single point repeated n times, which would make the density (8.4.4) uniform.) Then instead of the density (8.4.4), we define the *continuous inverse sensitivity mechanism* M_{cont} to have density

$$q(t \mid P_n) = \frac{\exp(-\frac{\varepsilon}{2} d_{f,\rho}(t, P_n))}{\int \exp(-\frac{\varepsilon}{2} d_{f,\rho}(s, P_n)) ds}. \quad (8.4.9)$$

While the parameter ρ adds complexity, setting it to be very small (say, $\rho = \frac{1}{n^2}$) is a reasonable practical default.

The continuous inverse sensitivity enjoys fairly strong error guarantees, as the next two propositions demonstrate, providing two prototypical results. (Exercises 8.11 and 8.12 show how to prove the propositions.) The first proposition shows that the inverse sensitivity mechanism is essentially never worse than the Laplace mechanism (8.1.3) when $\varepsilon \lesssim 1$.

Proposition 8.4.11. Let f be sample monotone (Definition 8.8) and have finite global sensitivity $\text{GS}(f) < \infty$. Then taking $\rho = 0$,

$$\mathbb{E} [|M_{\text{cont}}(P_n) - f(P_n)|] \leq \frac{1}{1 - e^{-\varepsilon/2}} \text{GS}(f).$$

As Example 8.1.3 shows, the standard Laplace mechanism M has error

$$\mathbb{E} [|M(P_n) - f(P_n)|] = \frac{\text{GS}(f)}{\varepsilon},$$

the same scaling Proposition 8.4.11 guarantees, because $1 - e^{-\varepsilon/2} = \varepsilon/2 + O(\varepsilon^2)$.

For the next proposition, which provides a more nuanced guarantee, we require local sensitivities for samples P'_n near P_n , and so we define the largest local sensitivity within Hamming distance K of the sample P_n by

$$L(K) := \sup_{P'_n \in \mathcal{P}_n} \{\text{LS}(f, P'_n) \mid d_{\text{ham}}(P_n, P'_n) \leq K\} = \sup_{P'_n \in \mathcal{P}_n} \{\omega_f(1, P'_n) \mid d_{\text{ham}}(P_n, P'_n) \leq K\},$$

where we recall the definition (8.4.2) of the local sensitivity of f . Then we have the following.

Proposition 8.4.12. *Let f be sample monotone (Definition 8.8) and have finite global sensitivity $\text{GS}(f) < \infty$. Then for any $\rho \geq 0$ and $K_n = \left\lceil \frac{4 \log(2n \text{GS}(f)/\rho)}{\varepsilon} \right\rceil$,*

$$\mathbb{E}[|M_{\text{cont}}(P_n) - f(P_n)|] \leq 2\rho + \frac{1}{1 - e^{-\varepsilon/2}} L(K_n).$$

Unpacking Proposition 8.4.12 a bit, let us make the default substitution $\rho = \frac{1}{n^2}$. Then because $1 - e^{-\varepsilon/2} = \varepsilon/2 + O(\varepsilon^2)$, for $\varepsilon \lesssim 1$ this yields

$$\mathbb{E}[|M_{\text{cont}}(P_n) - f(P_n)|] \lesssim \frac{1}{\varepsilon} \sup_{P'_n \in \mathcal{P}_n} \{\text{LS}(f, P'_n) \mid d_{\text{ham}}(P'_n, P_n) \leq K_n\} + \frac{1}{n^2},$$

where $K_n = \frac{4 \log \text{GS}(f) + 12 \log n}{\varepsilon} \lesssim \frac{1}{\varepsilon} \log n$ for large sample sizes n . Comparing this to the sketched lower bound (8.4.8), these quantities are of the same order whenever the moduli of continuity $\omega_f(k; P_n)$ are roughly additive and comparable near P_n , so that for $k \lesssim \frac{1}{\varepsilon}$ there is a chain $P_n^{(1)}, P_n^{(2)}, \dots, P_n^{(k)}$ with $d_{\text{ham}}(P_n^{(i)}, P_n^{(i+1)}) = 1$ and $\omega_f(k; P_n) \gtrsim \sum_{i=1}^k \text{LS}(f, P_n^{(i)})$ and $\text{LS}(f, P_n) \asymp \text{LS}(f, P'_n)$ for P'_n satisfying $d_{\text{ham}}(P_n, P'_n) \lesssim \frac{\log n}{\varepsilon}$. Under these conditions—which often require care to check, but which hold, for example, for mean estimation—we then obtain

$$\mathbb{E}[|M_{\text{cont}}(P_n) - f(P_n)|] \lesssim \omega_f\left(\frac{1}{\varepsilon}, P_n\right) + \frac{1}{n^2}.$$

8.5 Deferred proofs

8.5.1 Proof of Lemma 8.2.10

We prove the first statement of the lemma first. Let us assume there exists R such that $\|P - R\|_{\text{TV}} \leq \delta$ and $D_\infty(R\|Q) \leq \varepsilon$. Then for any set S we have

$$P(S) \leq R(S) + \delta \leq e^\varepsilon Q(S) + \delta, \quad \text{i.e.} \quad \log \frac{P(S) - \delta}{Q(S)} \leq \varepsilon,$$

which is equivalent to $D_\infty^\delta(P\|Q) \leq \varepsilon$. Now, let us assume that $D_\infty^\delta(P\|Q) \leq \varepsilon$, whence we must construct the distribution R .

We assume w.l.o.g. that P and Q have densities p, q , and define the sets

$$S := \{x : p(x) > e^\varepsilon q(x)\} \quad \text{and} \quad T := \{x : p(x) < q(x)\}.$$

On these sets, we have $0 \leq P(S) - e^\varepsilon Q(S) \leq \delta$ by assumption, and we then define a distribution R with density that we partially specify via

$$\begin{aligned} x \in S &\Rightarrow r(x) := e^\varepsilon q(x) < p(x) \\ x \in (T \cup S)^c &\Rightarrow r(x) := p(x) \leq e^\varepsilon q(x) \quad \text{and} \quad r(x) \geq q(x). \end{aligned}$$

Now, we note that $e^\varepsilon q(x) \geq p(x) \geq q(x)$ for $x \in (S \cup T)^c$, and thus

$$\begin{aligned} Q(S) + Q(S^c \cap T^c) &\leq e^\varepsilon Q(S) + P(S^c \cap T^c) \\ &= R(S) + R(S^c \cap T^c) \\ &= e^\varepsilon Q(S) + P(S^c \cap T^c) < P(S) + P(S^c \cap T^c). \end{aligned} \tag{8.5.1}$$

In particular, when $x \in T$, we may take the density r so that $p(x) \leq r(x) \leq q(x)$, as

$$R(S) + R(S^c \cap T^c) + P(T) < 1 \quad \text{and} \quad R(S) + R(S^c \cap T^c) + Q(T) > 1$$

by the inequalities (8.5.1), and so that $R(\mathcal{X}) = 1$. With this, we evidently have $r(x) \leq e^\varepsilon q(x)$ by construction, and because $S \subset T^c$, we have

$$R(T) - P(T) = P(T^c) - R(T^c) = P(S \cap T^c) - R(S \cap T^c) + P(S^c \cap T^c) - R(S^c \cap T^c) = P(S) - R(S),$$

where we have used that $r = p$ on $(T \cup S)^c$ by construction. Thus we find that

$$\begin{aligned} \|P - R\|_{\text{TV}} &= \frac{1}{2} \int_S |r - p| + \frac{1}{2} \int_T |r - p| = \frac{1}{2} (P(S) - R(S)) + \frac{1}{2} (R(T) - P(T)) \\ &= P(S) - R(S) = P(S) - e^\varepsilon Q(S) \leq \delta \end{aligned}$$

by assumption.

Now, we turn to the second statement of the lemma. We start with the easy direction, where we assume that P_0 and Q_0 satisfy $D_\infty(P_0 \| Q_0) \leq \varepsilon$ and $D_\infty(Q_0 \| P_0) \leq \varepsilon$ as well as $\|P - P_0\|_{\text{TV}} \leq \delta$ and $\|Q - Q_0\|_{\text{TV}} \leq \delta$. Then for any set S we have

$$P(S) \leq P_0(S) + \frac{\delta}{1 + e^\varepsilon} \leq e^\varepsilon Q_0(S) + \frac{\delta}{1 + e^\varepsilon} \leq e^\varepsilon Q(S) + e^\varepsilon \delta + \frac{\delta}{1 + e^\varepsilon},$$

or $D_\infty^\delta(P \| Q) \leq \varepsilon$. The other direction is similar.

We consider the converse direction, where we have both $D_\infty^\delta(P \| Q) \leq \varepsilon$ and $D_\infty^\delta(Q \| P) \leq \varepsilon$. Let us construct P_0 and Q_0 as in the statement of the lemma. Define the sets

$$S := \{x : p(x) > e^\varepsilon q(x)\} \quad \text{and} \quad S' := \{x : q(x) > e^\varepsilon p(x)\}$$

as well as the sets

$$T := \{x : e^\varepsilon q(x) \geq p(x) \geq q(x)\} \quad \text{and} \quad T' := \{x : e^{-\varepsilon} q(x) \leq p(x) < q(x)\},$$

so that S, S', T, T' are all disjoint, and $\mathcal{X} = S \cup S' \cup T \cup T'$. We begin by constructing intermediate measures—which end up not being probabilities— P_1 and Q_1 , which we modify slightly to actually construct P_0 and Q_0 . We first construct densities similar to our construction above for part (i), setting

$$\begin{aligned} x \in S &\Rightarrow p_1(x) := e^\varepsilon q_1(x), \quad q_1(x) := \frac{1}{1 + e^\varepsilon} (p(x) + q(x)) \\ x \in S' &\Rightarrow q_1(x) := e^\varepsilon p_1(x), \quad p_1(x) := \frac{1}{1 + e^\varepsilon} (p(x) + q(x)). \end{aligned}$$

Now, define the two quantities

$$\alpha := P(S) - P_1(S) = P(S) - \frac{e^\varepsilon}{1 + e^\varepsilon} (P(S) + Q(S)) = \frac{P(S) - e^\varepsilon Q(S)}{1 + e^\varepsilon} \leq \frac{\delta}{1 + e^\varepsilon}.$$

and similarly

$$\alpha' := Q(S') - Q_1(S') = \frac{Q(S') - e^\varepsilon P(S')}{1 + e^\varepsilon} \leq \frac{\delta}{1 + e^\varepsilon}.$$

Note also that we have $P(S) - P_1(S) = Q_1(S) - Q(S)$ and $Q(S') - Q_1(S') = P_1(S') - P(S')$ by construction.

We assume w.l.o.g. that $\alpha \geq \alpha'$, so that if $\beta = \alpha - \alpha' \geq 0$, we have $\beta \leq \frac{\delta}{1+e^\varepsilon}$, and we have the sandwiching

$$P_1(S) + P_1(S') + P(T \cup T') = P_1(S) + P_1(S') + 1 - P(S \cup S') = 1 - \beta < 1$$

because S and S' are disjoint and $T_{<} \cup T_{>} = (S \cup S')^c$, and similarly

$$Q_1(S) + Q_1(S') + Q(T \cup T') = Q_1(S) + Q_1(S') + 1 - Q(S \cup S') = 1 + \beta > 1.$$

Let $p_1 = p$ on the set $T \cup T'$ and similarly for $q_1 = q$. Then we have $P_1(\mathcal{X}) = 1 - \beta$, $Q_1(\mathcal{X}) = 1 + \beta$, and $|\log \frac{p_1}{q_1}| \leq \varepsilon$.

Now, note that $S \cup T = \{x : q_1(x) \geq p_1(x)\}$, and we have

$$\begin{aligned} Q_1(S) + Q_1(T) - P_1(S) - P_1(T) &= Q_1(S) + Q(T) - P_1(S) - P(T) \\ &\geq Q_1(S) + Q_1(S') + Q(T) + Q(T') - P_1(S) - P_1(S') - P(T) - P(T') = 2\beta. \end{aligned}$$

Now, (roughly) we decrease the density q_1 to q_0 on $S \cup T$ and increase p_1 to p_0 on $S \cup T$, while still satisfying $q_0 \geq p_0$ on $S \cup T$. In particular, we may choose the densities $q_0 = q_1$ on $T' \cup S'$ and $p_0 = p_1$ on $T' \cup S'$, while choosing q_0, p_0 so that

$$p_1(x) \leq p_0(x) \leq q_0(x) \leq q_1(x) \quad \text{on } S \cup T,$$

where

$$P_0(S \cup T) = P_1(S \cup T) + \beta \quad \text{and} \quad Q_0(S \cup T) = Q_1(S \cup T) - \beta. \quad (8.5.2)$$

With these choices, we evidently obtain $Q_0(\mathcal{X}) = P_0(\mathcal{X}) = 1$ and that $D_\infty(P_0 \| Q_0) \leq \varepsilon$ and $D_\infty(Q_0 \| P_0) \leq \varepsilon$ by construction. It remains to consider the variation distances. As $p_0 = p$ on T' , we have

$$\begin{aligned} \|P - P_0\|_{\text{TV}} &= \frac{1}{2} \int_S |p - p_0| + \frac{1}{2} \int_{S'} |p - p_0| + \frac{1}{2} \int_T |p - p_0| \\ &= \frac{1}{2} (P(S) - P_0(S)) + \frac{1}{2} (P_0(S') - P(S)) + \frac{1}{2} (P_0(T) - P(T)) \\ &\leq \frac{1}{2} \underbrace{(P(S) - P_1(S))}_{=\alpha} + \frac{1}{2} \underbrace{(P_0(S') - P(S))}_{=\alpha'} + \frac{1}{2} \underbrace{(P_0(T) - P(T))}_{\leq \beta}, \end{aligned}$$

where the $P_0(T) - P(T) \leq \beta$ claim follows because $p_1(x) = p(x)$ on T and by the increasing construction yielding equality (8.5.2), we have $P_0(T) - P(T) = P_0(T) - P_1(T) = \beta + P_1(S) - P_0(S) \leq \beta$. In particular, we have $\|P - P_0\|_{\text{TV}} \leq \frac{\alpha + \alpha'}{2} + \frac{\beta}{2} = \alpha \leq \frac{\delta}{1+e^\varepsilon}$. The argument that $\|Q - Q_0\|_{\text{TV}} \leq \frac{\delta}{1+e^\varepsilon}$ is similar.

8.6 Bibliography

Given the broad focus of this book, our treatment of privacy is necessarily somewhat brief, and there is substantial depth to the subject that we do not cover.

The initial development of randomized response began with Warner [190], who proposed randomized response in survey sampling as a way to collect sensitive data. This elegant idea remained in use for many years, and a generalization to data release mechanisms with bounded likelihood ratios—essentially, the local differential privacy definition 8.2—is due to Evfimievski et al. [86] in 2003 in the databases community. Dwork, McSherry, Nissim, and Smith [80] and the subsequent work of Dwork et al. [79] defined differential privacy and its (ϵ, δ) -approximate relaxation. A small industry of research has built out of these papers, with numerous extensions and developments.

Exponential mechanism is McSherry and Talwar [148].

The book of Dwork and Roth [78] surveys much of the field, from the perspective of computer science, as of 2014. Lemma 8.2.10 is due to Dwork et al. [81], and our proof is based on theirs.

8.7 Exercises

Exercise 8.1: Prove Proposition 8.2.1.

Exercise 8.2: Prove Proposition 8.4.7.

Exercise 8.3 (Laplace mechanisms versus randomized response): In this question, you will investigate using Laplace and randomized response mechanisms, as in Examples 8.1.3 and 8.1.1–8.1.2, to perform *locally* private estimation of a mean, and compare this with randomized-response based mechanisms.

We consider the following scenario: we have data $X_i \in [0, 1]$, drawn i.i.d., and wish to estimate the mean $\mathbb{E}[X]$ under local ϵ -differential privacy.

- (a) The Laplace mechanism simply sets $Z_i = X_i + W_i$ for $W_i \stackrel{\text{iid}}{\sim} \text{Laplace}(b)$ for some b . What choice of b guarantees ϵ -local differential privacy?
- (b) For your choice of b , let $\bar{Z}_n = \frac{1}{n} \sum_{i=1}^n Z_i$. Give $\mathbb{E}[(\bar{Z}_n - \mathbb{E}[X])^2]$.
- (c) A randomized response mechanism for this case is the following: first, we randomly round X_i to $\{0, 1\}$, by setting

$$\tilde{X}_i = \begin{cases} 1 & \text{with probability } X_i \\ 0 & \text{otherwise.} \end{cases}$$

Conditional on $\tilde{X}_i = x$, we then set

$$Z_i = \begin{cases} x & \text{with probability } \frac{e^\epsilon}{1+e^\epsilon} \\ 1-x & \text{with probability } \frac{1}{1+e^\epsilon}. \end{cases}$$

What is $\mathbb{E}[Z_i]$?

- (d) For the randomized response Z_i above, give constants a and b so that $aZ_i - b$ is unbiased for $\mathbb{E}[X]$, that is, $\mathbb{E}[aZ_i - b] = \mathbb{E}[X]$. Let $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n (aZ_i - b)$ be your mean estimator. What is $\mathbb{E}[(\hat{\theta}_n - \mathbb{E}[X])^2]$? Does this converge to the mean-square error of the sample mean $\mathbb{E}[(\bar{X}_n - \mathbb{E}[X])^2] = \text{Var}(X)/n$ as $\epsilon \uparrow \infty$?

(e) Now, it is time to compare the simple randomized response estimator from part (d) with the Laplace mechanism from part (b). For each of the following distributions, generate samples of size $N = 10, 100, 1000, 10000$, and then for $T = 25$ tests, compute the two estimators, both with $\varepsilon = 1$. Then plot the mean-squared error and confidence intervals for each of the two methods as well as the sample mean without any privacy.

- i. Uniform distribution: $X \sim \text{Uniform}[0, 1]$, with $\mathbb{E}[X] = 1/2$.
- ii. Bernoulli distribution: $X \sim \text{Bernoulli}(p)$, where $p = .1$.
- iii. Uniform distribution: $X \sim \text{Uniform}[,49, .51]$, with $\mathbb{E}[X] = 1/2$.

Do you prefer the Laplace or randomized response mechanism? In one sentence, why?

Exercise 8.4 (A more sophisticated randomized response scheme): Let us consider a more sophisticated randomized response scheme than that in Exercise 8.3. Define quantized values

$$b_0 = 0, b_1 = \frac{1}{k}, \dots, b_{k-1} = \frac{k-1}{k}, b_k = 1. \quad (8.7.1)$$

Now consider a randomized response estimator that, when $X \in [b_j, b_{j+1}]$ first rounds X randomly to $\tilde{X} \in \{b_j, b_{j+1}\}$ so that $\mathbb{E}[\tilde{X} | X] = X$. Conditional on $\tilde{X} = j$, we then set

$$Z = \begin{cases} j & \text{with probability } \frac{e^\varepsilon}{k+e^\varepsilon} \\ \text{Uniform}(\{0, \dots, k\} \setminus \{j\}) & \text{with probability } \frac{k}{k+e^\varepsilon}. \end{cases}$$

- (a) Give a and b so that $\mathbb{E}[aZ - b] = \mathbb{E}[X]$.
- (b) For your values of a and b above, let $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n (aZ_i - b)$. Give a (reasonably tight) bound on $\mathbb{E}[(\hat{\theta}_n - \mathbb{E}[X])^2]$.
- (c) For any given $\varepsilon > 0$, give (approximately) the k in the choice of the number of bins (8.7.1) that optimizes your bound, and (approximately) evaluate $\mathbb{E}[(\hat{\theta}_n - \mathbb{E}[X])^2]$ with your choice of k . As $\varepsilon \uparrow \infty$, does this converge to $\text{Var}(X)/n$?

Exercise 8.5 (Subsampling via divergence measures (Balle et al. [14])): The *hockey stick* divergence functional, defined for $\alpha \geq 1$, is $\phi_\alpha(t) = [1 - \alpha t]_+$. It is straightforward to relate this to (ε, δ) -differential privacy via Definition 8.6: two distributions P and Q are (ε, δ) -close if and only if their ϕ_{e^ε} -divergences are less than δ , i.e., if and only if

$$D_{\phi_{e^\varepsilon}}(P \| Q) \leq \delta \quad \text{and} \quad D_{\phi_{e^\varepsilon}}(Q \| P) \leq \delta.$$

(In your answer to this question, feel free to use $D_\alpha(P \| Q)$ as a shorthand for $D_{\phi_\alpha}(P \| Q)$.)

- (a) Let P_0, P_1, Q_1 be any three distributions, and for some $q \in [0, 1]$ and $\alpha \geq 1$, define $P = (1 - q)P_0 + qP_1$ and $Q = (1 - q)P_0 + qQ_1$. Let $\alpha' = 1 + q(\alpha - 1) = (1 - q) + q\alpha$ and $\theta = \alpha'/\alpha \leq 1$. Show that

$$D_{\phi_{\alpha'}}(P \| Q) = qD_{\phi_\alpha}((1 - \theta)P_0 + \theta P_1 \| Q_1).$$

- (b) Let $\varepsilon > 0$ and define $\varepsilon(q) = \log(1 + q(e^\varepsilon - 1))$. Show that

$$D_{\phi_{e^{\varepsilon(q)}}}(P \| Q) \leq q \max \{D_{\phi_{e^\varepsilon}}(P_0 \| Q_1), D_{\phi_{e^\varepsilon}}(P_1 \| Q_1)\}.$$

Exercise 8.6 (Subsampling and privacy amplification (Balle et al. [14])): Consider the following subsampling approach to privacy. Assume that we have a private (randomized) algorithm, represented by \mathcal{A} , that acts on samples of size m and guarantees (ε, δ) -differential privacy. The subsampling mechanism is then defined as follows: given a sample X_1^n of size $n > m$, choose a subsample X_{sub} of size m uniformly at random from X_1^n , and then release $Z = \mathcal{A}(X_{\text{sub}})$.

- (a) Use the results of parts (a) and (b) in Exercise 8.5 to show that Z is $(\varepsilon(q), \delta q)$ -differentially private, where $q = m/n$ and $\varepsilon(q) = \log(1 + q(e^\varepsilon - 1))$.
- (b) Show that if $\varepsilon \leq 1$, then Z is $((e - 1)q\varepsilon, q\delta)$ -differentially private, and if $\varepsilon \leq \frac{1}{2}$, then Z is $(2(\sqrt{e} - 1)q\varepsilon, q\delta)$ -differentially private. *Hint:* Argue that for any $T > 0$, one has $e^t - 1 \leq (e^T - 1)^{\frac{t}{T}}$ for all $t \in [0, T]$.

Exercise 8.7 (Concentration and privacy composition): In this question, we give an alternative to the privacy composition approaches we exploit in Section 8.3.2. Consider an identical scenario to that in Fig. 8.1, and begin by assuming that each channel Q_i is ε -differentially private with density q_i , and let $Q^{(b)}$ be shorthand for $Q(\cdot \mid x^{(b)})$. Define the log-likelihood ratio

$$L^{(b)}(Z_1^k) := \sum_{i=1}^k \log \frac{q_i^{(b)}(Z_i)}{q_i^{(1-b)}(Z_i)}.$$

- (a) Let P, Q be any two distributions satisfying $D_\infty(P\|Q) \leq \varepsilon$ and $D_\infty(Q\|P) \leq \varepsilon$, i.e., that $\log \frac{P(A)}{Q(A)} \in [-\varepsilon, \varepsilon]$ for all sets A . Show that

$$D_{\text{kl}}(P\|Q) \leq \varepsilon(e^\varepsilon - 1).$$

- (b) Let $Q^{(b)}$ denote the joint distribution of Z_1, \dots, Z_k when bit b holds in the privacy game in Fig. 8.1. Show that

$$\mathbb{E}_b[L^{(b)}(Z_1^k)] \leq k\varepsilon(e^\varepsilon - 1)$$

where \mathbb{E}_b denotes expectation under $Q^{(b)}$, and that for all $t \geq 0$,

$$Q^{(b)}\left(L^{(b)}(Z_1^k) \geq k\varepsilon(e^\varepsilon - 1) + t\right) \leq \exp\left(-\frac{t^2}{2k\varepsilon^2}\right).$$

Conclude that for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over $Z_1^k \sim Q^{(b)}$,

$$L^{(b)}(Z_1^k) \leq k(e^\varepsilon - 1)\varepsilon + \sqrt{2k \log \frac{1}{\delta}} \cdot \varepsilon.$$

- (c) Argue that for any (measurable) set A ,

$$Q^{(b)}(Z_1^k \in A) \leq e^{\varepsilon(k, \delta)} \cdot Q^{(1-b)}(Z_1^k \in A) + \delta$$

for all $\delta \in [0, 1]$, where $\varepsilon(k, \delta) = k\varepsilon(e^\varepsilon - 1) + \sqrt{2k \log \frac{1}{\delta}} \cdot \varepsilon$.

- (d) Conclude the following tighter variant of Corollary 8.3.3: if each channel in Fig. 8.1 is ε -differentially private, then the composition of k such channels is

$$\left(k\varepsilon(e^\varepsilon - 1) + \sqrt{2k \log \frac{1}{\delta}} \cdot \varepsilon, \delta\right)$$

differentially private for all $\delta > 0$.

As an aside, a completely similar derivation yields the following tighter analogue of Theorem 8.3.4: if each channel is (ε, δ) -differentially private, then their composition is

$$\left(k\varepsilon(e^\varepsilon - 1) + \sqrt{2k \log \frac{1}{\delta_0}} \cdot \varepsilon, \delta_0 + \frac{k\delta}{1 + e^\varepsilon} \right)$$

differentially private for all $\delta_0 > 0$.

Exercise 8.8 (One-dimensional minimization with inverse sensitivity): Consider the private minimization of the one dimensional loss $\ell(\theta, x)$ (for $\theta \in \Theta \subset \mathbb{R}$), where we wish to estimate

$$\hat{\theta}(P_n) \in \operatorname{argmin}_{\theta} \{P_n \ell(\theta, X) := \frac{1}{n} \sum_{i=1}^n \ell(\theta, X_i)\},$$

where we recall the notation from Chapters 4 and 6. Assume that the loss ℓ is convex, differentiable in θ , and that it satisfies the Lipschitz-type guarantees that there exist constants $0 < L_0 \leq L_1 < \infty$

$$[-L_0, L_0] \subset \{\ell'(\theta, x)\}_{x \in \mathcal{X}} \subset [-L_1, L_1] \quad (8.7.2)$$

for all $\theta \in \Theta$ and that $\{\ell'(\theta, x)\}_{x \in \mathcal{X}}$ is an interval. (That is, the set of potential derivatives $\ell'(\theta, x)$ as x varies includes $[-L_0, L_0]$, is convex, and $|\ell'(\theta, x)| \leq L_1$ for all $\theta \in \Theta, x \in \mathcal{X}$.)

(a) Let the loss ℓ be the Huber loss $\ell(\theta, x) = h_u(\theta - x)$ for some fixed $u > 0$, where

$$h_u(t) = \begin{cases} \frac{1}{2u} t^2 & \text{if } |t| \leq u \\ |t| + \frac{u}{2} & \text{if } |t| \geq u. \end{cases}$$

When $\mathcal{X} = \mathbb{R}$, show that ℓ satisfies the containment (8.7.2) with $L_0 = L_1 = 1$.

(b) Let the loss ℓ be the absolute value $\ell(\theta, x) = |\theta - x|$, where we abuse notation to call $\{\ell'(\theta, x)\}_{x=\theta} = [-1, 1]$ (the subdifferential). When $\mathcal{X} = \mathbb{R}$, show that ℓ satisfies the containment (8.7.2) with $L_0 = L_1 = 1$.

(c) Let $d_{\hat{\theta}}$ be the inverse sensitivity (8.4.3) for the minimizer $\hat{\theta}(P_n)$, which is the solution (in θ) to $P_n \ell'(\theta, X) = 0$. Assuming inequality (8.7.2) holds, show that

$$\left\lceil \frac{n|P_n \ell'(\theta, X)|}{2L_1} \right\rceil \leq d_{\hat{\theta}}(\theta, P_n) \leq \left\lceil \frac{n|P_n \ell'(\theta, X)|}{L_0} \right\rceil.$$

(d) Show that the function

$$d(\theta, P_n) := \left\lceil \frac{n|P_n \ell'(\theta, X)|}{2L_1} \right\rceil$$

is 1-Lipschitz with respect to the Hamming metric in P_n .

The Lipschitz behavior of $d(\theta, P_n)$ in part (d) makes this a computationally attractive alternative to the pure inverse sensitivity (8.4.3) and associated mechanism with density (8.4.4).

Exercise 8.9 (Estimating means with inverse sensitivity mechanisms): In this question, we compare behavior of mean estimation under differential privacy with the Laplace mechanism and the inverse sensitivity-type mechanism in Example 8.4.8. Let $\mathcal{X} = [-1, 1]$ be the data space and consider estimating the mean \bar{x}_n of $x_1^n \in \mathcal{X}^n$.

- (a) Implement the Laplace mechanism (8.1.3) for this problem. Fix $n = 200$ and repeat the following experiment 50 times. For $\varepsilon = .1, .5, 1, 2$, generate a sample $x_1^n \in \mathcal{X}^n$ (from whatever distribution you like), then estimate \bar{x}_n using the Laplace mechanism. Give a table of the mean squared errors $(\bar{x}_n - M(x_1^n))^2$.
- (b) Implement the inverse sensitivity mechanism using the approximation in Example 8.4.8. Repeat the experiment in part (a).
- (c) Compare the results.

Exercise 8.10 (Estimating medians with the inverse sensitivity mechanism): The data at <https://stats311.stanford.edu/data/salaries.txt> contains approximately 250,000 salaries from the University of California Schools between 2011 and 2014. Assuming that the maximum salary is $3 \cdot 10^6$ and minimum is 0 (so the data $x \in [0, 3 \cdot 10^6]$), implement the inverse sensitivity mechanism for the median as in Example 8.4.9. Repeat the following 20 times: for each of $\varepsilon = .0625, .125, .25, .5, 1, 2$, estimate the median using the inverse sensitivity mechanism with ε -differential privacy. Compute the mean absolute errors across the 20 experiments for each ε .

Exercise 8.11 (Shells and accuracy in inverse sensitivity): Let $f : \mathcal{P}_n \rightarrow \mathbb{R}$ be sample monotone (Def. 8.8) and $\rho \geq 0$. Let $M = M_{\text{cont}}$ be the continuous inverse sensitivity mechanism with density (8.4.9). Define the upper and lower shells

$$S_{k+} = \{t > f(P_n) \mid d_{f,\rho}(t, P_n) = k\} \quad \text{and} \quad S_{k-} = \{t < f(P_n) \mid d_{f,\rho}(t, P_n) = k\},$$

and the upper and lower moduli of continuity (values in the shells $S_{k\pm}$)

$$\omega^+(k) := \sup\{t \in S_{k+}\} - f(P_n) \quad \text{and} \quad \omega^-(k) := f(P_n) - \inf\{t \in S_{k-}\}.$$

Let $S_0 = \{t \in \mathbb{R} \mid |f(P_n) - t| \leq \rho\}$.

- (a) Justify the inequality

$$\begin{aligned} & \mathbb{E}[|M(P_n) - f(P_n)|] \\ & \leq \mathbb{P}(M(P_n) \in S_0)\rho + \sum_{k=1}^n \mathbb{P}(M(P_n) \in S_{k+})(\omega^+(k) + \rho) + \sum_{k=1}^n \mathbb{P}(M(P_n) \in S_{k-})(\omega^-(k) + \rho). \end{aligned}$$

- (b) Bound $\mathbb{P}(M(P_n) \in S_{k+})$ and $\mathbb{P}(M(P_n) \in S_{k-})$, and using these bounds demonstrate that

$$\begin{aligned} & \mathbb{E}[|M(P_n) - f(P_n)|] \\ & \leq \rho + \frac{\sum_{k=1}^n \omega^+(k) \cdot (\omega^+(k) - \omega^+(k-1))e^{-k\varepsilon/2}}{\rho + \sum_{k=1}^n (\omega^+(k) - \omega^+(k-1))e^{-k\varepsilon/2} + \sum_{k=1}^n (\omega^-(k) - \omega^-(k-1))e^{-k\varepsilon/2}} \\ & \quad + \frac{\sum_{k=1}^n \omega^-(k) \cdot (\omega^-(k) - \omega^-(k-1))e^{-k\varepsilon/2}}{\rho + \sum_{k=1}^n (\omega^-(k) - \omega^-(k-1))e^{-k\varepsilon/2} + \sum_{k=1}^n (\omega^+(k) - \omega^+(k-1))e^{-k\varepsilon/2}} \end{aligned}$$

- (c) Show that

$$\sum_{k=1}^n [(\omega^-(k) - \omega^-(k-1)) + \omega^+(k) - \omega^+(k-1)] e^{-k\varepsilon/2} \geq (1 - e^{-\varepsilon/2}) \sum_{k=1}^n (\omega^+(k) + \omega^-(k)) e^{-k\varepsilon/2}.$$

Exercise 8.12 (Accuracy of the inverse sensitivity mechanism): In this question, we prove Propositions 8.4.11 and 8.4.12. Let the conditions and notation of Exercise 8.11 hold. Recall the definition

$$L(K) := \sup_{P'_n \in \mathcal{P}_n} \{ \text{LS}(f, P'_n) \mid d_{\text{ham}}(P_n, P'_n) \leq K \}.$$

(a) Use Exercise 8.11.(b) and (c) to show that for any $K \in \mathbb{N}$,

$$\begin{aligned} \mathbb{E} [|M_{\text{cont}}(P_n) - f(P_n)|] &\leq \rho + \frac{L(K)}{1 - e^{-\varepsilon/2}} \cdot \frac{\sum_{k=1}^K (\omega^+(k) + \omega^-(k)) e^{-k\varepsilon/2}}{\sum_{k=1}^n (\omega^+(k) + \omega^-(k)) e^{-k\varepsilon/2}} \\ &\quad + \frac{\text{GS}(f)}{\rho} \sum_{k=K+1}^n (\omega^+(k) + \omega^-(k)) e^{-k\varepsilon/2}. \end{aligned}$$

(b) Choose values for ρ and K to show that $\mathbb{E} [|M_{\text{cont}}(P_n) - f(P_n)|] \leq \frac{1}{1 - e^{-\varepsilon/2}} \text{GS}(f)$, giving Proposition 8.4.11.

(c) Prove Proposition 8.4.12.

Exercise 8.13 (Subsampling and Rényi privacy): We would like to estimate the mean $\mathbb{E}[X]$ of $X \sim P$, where $X \in \mathbb{B} = \{x \in \mathbb{R}^d \mid \|x\|_2 \leq 1\}$, the ℓ_2 -ball in \mathbb{R}^d . We investigate the extent to which subsampling of a dataset can *improve* privacy by providing some additional anonymity. Consider the following mechanism for estimating (scaled) multiples of this mean: for a dataset $\{X_1, \dots, X_n\}$, we let $S_i \in \{0, 1\}$ be i.i.d. Bernoulli(q), that is, $\mathbb{E}[S_i] = q$, and then consider the algorithm

$$Z = \sum_{i=1}^n X_i S_i + \sigma W, \quad W \sim \mathcal{N}(0, I_d). \quad (8.7.3)$$

In this question, we investigate the Rényi privacy properties of the subsampling (8.7.3). (Recall the Rényi divergence of Definition 8.4, $D_\alpha(P \| Q) = \frac{1}{\alpha-1} \log \int (p/q)^\alpha q$.)

We consider a slight variant of Rényi privacy, where we define data matrices X and X' to be adjacent if $X \in \mathbb{R}^{d \times n}$ and $X' \in \mathbb{R}^{d \times n-1}$ where X' is X with a single column removed. Then a mechanism is (ε, α) -Rényi private against single removals if and only if

$$D_\alpha(Q(\cdot \mid X) \| Q(\cdot \mid X')) \leq \varepsilon \quad \text{and} \quad D_\alpha(Q(\cdot \mid X') \| Q(\cdot \mid X)) \leq \varepsilon \quad (8.7.4)$$

for all neighboring X and X' consisting of samples of size n and $n-1$, respectively.

(a) Let $Q(\cdot \mid X)$ and $Q(\cdot \mid X')$ denote the channels for the mechanism (8.7.3) with data matrices $X = [x_1 \cdots x_{n-1} \ x]$ and $X' = [x_1 \cdots x_{n-1}] \in \mathbb{R}^{d \times n}$. Let P_μ denote the normal distribution $\mathcal{N}(\mu, \sigma^2 I)$ with mean μ and covariance $\sigma^2 I$ on \mathbb{R}^d . Show that for any $\alpha \in (1, \infty)$,

$$D_\alpha(Q(\cdot \mid X) \| Q(\cdot \mid X')) \leq D_\alpha(qP_x + (1-q)P_0 \| P_0)$$

and

$$D_\alpha(Q(\cdot \mid X') \| Q(\cdot \mid X)) \leq D_\alpha(P_0 \| qP_x + (1-q)P_0).$$

(b) Show that for the Rényi $\alpha = 2$ -divergence,

$$D_2(qP_x + (1-q)P_0 \| P_0) \leq \log \left(1 + q^2 \left(\exp(\|x\|_2^2 / \sigma^2) - 1 \right) \right) \quad \text{and}$$

$$D_2(P_0 \| qP_x + (1-q)P_0) \leq \log \left(1 + \frac{q^2}{1-q} \left(\exp(\|x\|_2^2 / \sigma^2) - 1 \right) \right).$$

(Hint: Example 8.2.2.)

Consider two mechanisms for computing a sample mean \bar{X}_n of vectors, where $\|x_i\|_2 \leq b$ for all i . The first is to repeat the following T times: for $t = 1, 2, \dots, T$,

- i. Draw $S \in \{0, 1\}^n$ with $S_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(q)$
- ii. Set $Z_t = \frac{1}{nq}(XS + \sigma_{\text{sub}}W_t)$, where $W_t \stackrel{\text{iid}}{\sim} \mathbf{N}(0, I)$, as in (8.7.3).

Then set $Z_{\text{sub}} = \frac{1}{T} \sum_{t=1}^T Z_t$. The other mechanism is to simply set $Z_{\text{Gauss}} = \bar{X}_n + \sigma_{\text{Gauss}}W$ for $W \sim \mathbf{N}(0, I)$.

- (c) What level of privacy does Z_{sub} have? That is, Z_{sub} is $(\varepsilon, 2)$ -Rényi private (against single removals (8.7.4)). Give a tight upper bound on ε .
- (d) What level of $(\varepsilon, 2)$ -Rényi privacy does Z_{Gauss} provide?
- (e) Fix $\varepsilon > 0$, and assume that each mechanism Z_{sub} and Z_{Gauss} have parameters chosen so that they are $(\varepsilon, 2)$ -Rényi private. Optimize over $T, q, n, \sigma_{\text{sub}}$ in the subsampling mechanism and σ_{Gauss} in the Gaussian mechanism, and provide the sharpest bound you can on

$$\mathbb{E}[\|Z_{\text{sub}} - \bar{X}_n\|_2^2] \quad \text{and} \quad \mathbb{E}[\|Z_{\text{Gauss}} - \bar{X}_n\|_2^2].$$

You may assume $\|x_i\|_2 = b$ for all i . (In your derivation, to avoid annoying constants, you should replace $\log(1+t)$ with its upper bound, $\log(1+t) \leq t$, which is fairly sharp for $t \approx 0$.)

Part II

Fundamental limits and optimality

JCD Comment: Put a brief commentary here. Some highlights:

- i. Minimax lower bounds (both local and global) using Le Cam's, Fano's, and Assouad's methods. Worked out long example with nonparametric regression.
- ii. Strong data processing inequalities, along with some bounds on them (constrained risk inequalities).
- iii. Functionals for lower bounds perhaps

Chapter 9

Minimax lower bounds: the Le Cam, Fano, and Assouad methods

Understanding the fundamental limits of estimation and optimization procedures is important for a multitude of reasons. Indeed, developing bounds on the performance of procedures can give complementary insights. By exhibiting fundamental limits of performance (perhaps over restricted classes of estimators), it is possible to guarantee that an algorithm we have developed is optimal, so that searching for estimators with better statistical performance will have limited returns, though searching for estimators with better performance in other metrics may be interesting. Moreover, exhibiting refined lower bounds on the performance of estimators can also suggest avenues for developing alternative, new optimal estimators; lower bounds need not be a fully pessimistic exercise.

In this chapter, we define and then discuss techniques for lower-bounding the minimax risk, giving three standard techniques for deriving minimax lower bounds that have proven fruitful in a variety of estimation problems [194]. In addition to reviewing these standard techniques—the Le Cam, Fano, and Assouad methods—we present a few simplifications and extensions that may make them more “user friendly.” Finally, the concluding sections of the chapter (Sections 10.1 and 10.2) present extensions of the ideas to *nonparametric problems*, where the effective number of parameters to estimate grows with the sample size n ; this culminates with an essentially geometric treatment of information and divergence measures directly relating covering and packing numbers to estimation.

9.1 Basic framework and minimax risk

Our first step here is to establish the minimax framework we use. When we study classical estimation problems, we use a standard version of minimax risk; we will also show how minimax bounds can be used to study optimization problems, in which case we use a specialization of the general minimax risk that we call minimax *excess* risk (while minimax risk handles this case, it is important enough that we define additional notation).

Let us begin by defining the standard minimax risk, deferring temporarily our discussion of minimax excess risk. Throughout, we let \mathcal{P} denote a class of distributions on a sample space \mathcal{X} , and let $\theta : \mathcal{P} \rightarrow \Theta$ denote a function defined on \mathcal{P} , that is, a mapping $P \mapsto \theta(P)$. The goal is to estimate the parameter $\theta(P)$ based on observations X_i drawn from the (unknown) distribution P . In certain cases, the parameter $\theta(P)$ uniquely determines the underlying distribution; for example, if we attempt to estimate a normal mean θ from the family $\mathcal{P} = \{\mathcal{N}(\theta, \sigma^2) : \theta \in \mathbb{R}\}$ with

known variance σ^2 , then $\theta(P) = \mathbb{E}_P[X]$ uniquely determines distributions in \mathcal{P} . In other scenarios, however, θ does not uniquely determine the distribution: for instance, we may be given a class of densities \mathcal{P} on the unit interval $[0, 1]$, and we wish to estimate $\theta(P) = \int_0^1 (p'(t))^2 dt$, where p is the density of P . Such problems arise, for example, in estimating the uniformity of the distribution of a species over an area (large $\theta(P)$ indicates an irregular distribution). In this case, θ does not parameterize P , so we take a slightly broader viewpoint of estimating functions of distributions in these notes.

The space Θ in which the parameter $\theta(P)$ takes values depends on the underlying statistical problem; as an example, if the goal is to estimate the univariate mean $\theta(P) = \mathbb{E}_P[X]$, we have $\Theta \subset \mathbb{R}$. To evaluate the quality of an estimator $\hat{\theta}$, we let $\rho : \Theta \times \Theta \rightarrow \mathbb{R}_+$ denote a (semi)metric on the space Θ , which we use to measure the error of an estimator for the parameter θ , and let $\Phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ be a non-decreasing function with $\Phi(0) = 0$ (for example, $\Phi(t) = t^2$).

For a distribution $P \in \mathcal{P}$, we assume we receive i.i.d. observations X_i drawn according to some P , and based on these $\{X_i\}$, the goal is to estimate the unknown parameter $\theta(P) \in \Theta$. For a given estimator $\hat{\theta}$ —a measurable function $\hat{\theta} : \mathcal{X}^n \rightarrow \Theta$ —we assess the quality of the estimate $\hat{\theta}(X_1, \dots, X_n)$ in terms of the risk

$$\mathbb{E}_P \left[\Phi(\rho(\hat{\theta}(X_1, \dots, X_n), \theta(P))) \right].$$

For instance, for a univariate mean problem with $\rho(\theta, \theta') = |\theta - \theta'|$ and $\Phi(t) = t^2$, this risk is the mean-squared error. As the distribution P is varied, we obtain the *risk functional* for the problem, which gives the risk of any estimator $\hat{\theta}$ for the family \mathcal{P} .

For any fixed distribution P , there is always a trivial estimator of $\theta(P)$: simply return $\theta(P)$, which will have minimal risk. Of course, this “estimator” is unlikely to be good in any real sense, and it is thus important to consider the risk functional not in a pointwise sense (as a function of individual P) but to take a more global view. One approach to this is Bayesian: we place a prior π on the set of possible distributions \mathcal{P} , viewing $\theta(P)$ as a random variable, and evaluate the risk of an estimator $\hat{\theta}$ taken in expectation with respect to this prior on P . Another approach, first suggested by Wald [189], which is to choose the estimator $\hat{\theta}$ minimizing the maximum risk

$$\sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\Phi(\rho(\hat{\theta}(X_1, \dots, X_n), \theta(P))) \right].$$

An optimal estimator for this metric then gives the *minimax risk*, which is defined as

$$\mathfrak{M}_n(\theta(\mathcal{P}), \Phi \circ \rho) := \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\Phi(\rho(\hat{\theta}(X_1, \dots, X_n), \theta(P))) \right], \quad (9.1.1)$$

where we take the supremum (worst-case) over distributions $P \in \mathcal{P}$, and the infimum is taken over all estimators $\hat{\theta}$. Here the notation $\theta(\mathcal{P})$ indicates that we consider parameters $\theta(P)$ for $P \in \mathcal{P}$ and distributions in \mathcal{P} .

In some scenarios, we study a specialized notion of risk appropriate for optimization problems (and statistical problems in which all we care about is prediction). In these settings, we assume there exists some loss function $\ell : \Theta \times \mathcal{X} \rightarrow \mathbb{R}$, where for an observation $x \in \mathcal{X}$, the value $\ell(\theta; x)$ measures the instantaneous loss associated with using θ as a predictor. In this case, we define the risk

$$L_P(\theta) := \mathbb{E}_P[\ell(\theta; X)] = \int_{\mathcal{X}} \ell(\theta; x) dP(x) \quad (9.1.2)$$

as the expected loss of the vector θ . (See, e.g., Chapter 5 of the lectures by Shapiro, Dentcheva, and Ruszczyński [168], or work on stochastic approximation by Nemirovski et al. [150].)

Example 9.1.1 (Support vector machines): In linear classification problems, we observe pairs $z = (x, y)$, where $y \in \{-1, 1\}$ and $x \in \mathbb{R}^d$, and the goal is to find a parameter $\theta \in \mathbb{R}^d$ so that $\text{sign}(\langle \theta, x \rangle) = y$. A convex loss surrogate for this problem is the hinge loss $\ell(\theta; z) = [1 - y\langle \theta, x \rangle]_+$; minimizing the associated risk functional (9.1.2) over a set $\Theta = \{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq r\}$ gives the support vector machine [55]. \diamond

Example 9.1.2 (Two-stage stochastic programming): In operations research, one often wishes to allocate resources to a set of locations $\{1, \dots, m\}$ before seeing demand for the resources. Suppose that the (unobserved) sample x consists of the pair $x = (C, v)$, where $C \in \mathbb{R}^{m \times m}$ corresponds to the prices of shipping a unit of material, so $c_{ij} \geq 0$ gives the cost of shipping from location i to j , and $v \in \mathbb{R}^m$ denotes the value (price paid for the good) at each location. Letting $\theta \in \mathbb{R}_+^m$ denote the amount of resources allocated to each location, we formulate the loss as

$$\ell(\theta; x) := \inf_{r \in \mathbb{R}^m, T \in \mathbb{R}^{m \times m}} \left\{ \sum_{i,j} c_{ij} T_{ij} - \sum_{i=1}^m v_i r_i \mid r_i = \theta_i + \sum_{j=1}^m T_{ji} - \sum_{j=1}^m T_{ij}, T_{ij} \geq 0, \sum_{j=1}^m T_{ij} \leq \theta_i \right\}.$$

Here the variables T correspond to the goods transported to and from each location (so T_{ij} is goods shipped from i to j), and we wish to minimize the cost of our shipping and maximize the profit. By minimizing the risk (9.1.2) over a set $\Theta = \{\theta \in \mathbb{R}_+^m : \sum_i \theta_i \leq b\}$, we maximize our expected reward given a budget constraint b on the amount of allocated resources. \diamond

For a (potentially random) estimator $\hat{\theta} : \mathcal{X}^n \rightarrow \Theta$ given access to a sample X_1, \dots, X_n , we may define the associated maximum *excess risk* for the family \mathcal{P} by

$$\sup_{P \in \mathcal{P}} \left\{ \mathbb{E}_P \left[L_P(\hat{\theta}(X_1, \dots, X_n)) \right] - \inf_{\theta \in \Theta} L(\theta) \right\},$$

where the expectation is taken over X_i and any randomness in the procedure $\hat{\theta}$. This expression captures the difference between the (expected) risk performance of the procedure $\hat{\theta}$ and the best possible risk, available if the distribution P were known ahead of time. The *minimax excess risk*, defined with respect to the loss ℓ , domain Θ , and family \mathcal{P} of distributions, is then defined by the best possible maximum excess risk,

$$\mathfrak{M}_n(\Theta, \mathcal{P}, \ell) := \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \left\{ \mathbb{E}_P \left[L_P(\hat{\theta}(X_1, \dots, X_n)) \right] - \inf_{\theta \in \Theta} L_P(\theta) \right\}, \quad (9.1.3)$$

where the infimum is taken over all estimators $\hat{\theta} : \mathcal{X}^n \rightarrow \Theta$ and the risk L_P is implicitly defined in terms of the loss ℓ . The techniques for providing lower bounds for the minimax risk (9.1.1) or the excess risk (9.1.3) are essentially identical; we focus for the remainder of this section on techniques for providing lower bounds on the minimax risk.

9.2 Preliminaries on methods for lower bounds

There are a variety of techniques for providing lower bounds on the minimax risk (9.1.1). Each of them transforms the maximum risk by lower bounding it via a Bayesian problem (e.g. [115, 134, 136]), then proving a lower bound on the performance of all possible estimators for the Bayesian problem (it is often the case that the worst case Bayesian problem is equivalent to the original

minimax problem [134]). In particular, let $\{P_v\} \subset \mathcal{P}$ be a collection of distributions in \mathcal{P} indexed by v and π be any probability mass function over v . Then for any estimator $\hat{\theta}$, the maximum risk has lower bound

$$\sup_{P \in \mathcal{P}} \mathbb{E}_P [\Phi(\rho(\hat{\theta}(X_1^n), \theta(P)))] \geq \sum_v \pi(v) \mathbb{E}_{P_v} [\Phi(\rho(\hat{\theta}(X_1^n), \theta(P_v)))] .$$

While trivial, this lower bound serves as the departure point for each of the subsequent techniques for lower bounding the minimax risk.

9.2.1 From estimation to testing

A standard first step in proving minimax bounds is to “reduce” the estimation problem to a testing problem [194, 192, 182]. The idea is to show that the probability of error in testing problems lower bounds the estimation risk, and we can develop tools for the former. We use two types of testing problems: one a multiple hypothesis test and the second based on multiple binary hypothesis tests.

Given an index set \mathcal{V} of finite cardinality, consider a family of distributions $\{P_v\}_{v \in \mathcal{V}}$ contained within \mathcal{P} . This family induces a collection of parameters $\{\theta(P_v)\}_{v \in \mathcal{V}}$; we call the family a 2δ -packing in the ρ -semimetric if

$$\rho(\theta(P_v), \theta(P_{v'})) \geq 2\delta \quad \text{for all } v \neq v'.$$

We use this family to define the *canonical hypothesis testing problem*:

- first, nature chooses V according to the uniform distribution over \mathcal{V} ;
- second, conditioned on the choice $V = v$, the random sample $X = X_1^n = (X_1, \dots, X_n)$ is drawn from the n -fold product distribution P_v^n .

Given the observed sample X , the goal is to determine the value of the underlying index v . We refer to any measurable mapping $\Psi : \mathcal{X}^n \rightarrow \mathcal{V}$ as a test function. Its associated error probability is $\mathbb{P}(\Psi(X_1^n) \neq V)$, where \mathbb{P} denotes the joint distribution over the random index V and X . In particular, if we set $\bar{P} = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P_v$ to be the mixture distribution, then the sample X is drawn (marginally) from \bar{P} , and our hypothesis testing problem is to determine the randomly chosen index V given a sample from this mixture \bar{P} .

With this setup, we obtain the classical reduction from estimation to testing.

Proposition 9.2.1. *The minimax error (9.1.1) has lower bound*

$$\mathfrak{M}_n(\theta(\mathcal{P}), \Phi \circ \rho) \geq \Phi(\delta) \inf_{\Psi} \mathbb{P}(\Psi(X_1, \dots, X_n) \neq V), \quad (9.2.1)$$

where the infimum ranges over all testing functions.

Proof To see this result, fix an arbitrary estimator $\hat{\theta}$. Suppressing dependence on X throughout the derivation, first note that it is clear that for any fixed θ , we have

$$\mathbb{E}[\Phi(\rho(\hat{\theta}, \theta))] \geq \mathbb{E}[\Phi(\delta) \mathbf{1}\{\rho(\hat{\theta}, \theta) \geq \delta\}] = \Phi(\delta) \mathbb{P}(\rho(\hat{\theta}, \theta) \geq \delta),$$

where the final inequality follows because Φ is non-decreasing. Now, let us define $\theta_v = \theta(P_v)$, so that $\rho(\theta_v, \theta_{v'}) \geq 2\delta$ for $v \neq v'$. By defining the testing function

$$\Psi(\hat{\theta}) := \operatorname{argmin}_{v \in \mathcal{V}} \{\rho(\hat{\theta}, \theta_v)\},$$

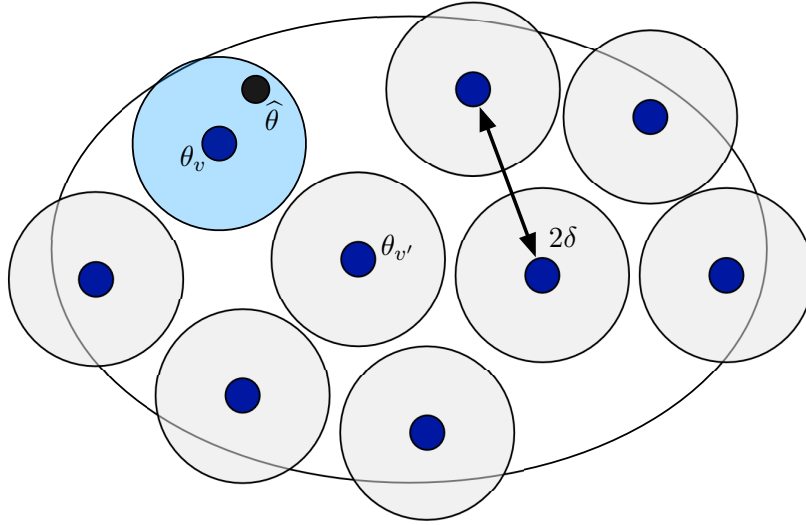


Figure 9.1. Example of a 2δ -packing of a set. The estimate $\hat{\theta}$ is contained in at most one of the δ -balls around the points θ_v .

breaking ties arbitrarily, we have that $\rho(\hat{\theta}, \theta_v) < \delta$ implies that $\Psi(\hat{\theta}) = v$ because of the triangle inequality and 2δ -separation of the set $\{\theta_v\}_{v \in \mathcal{V}}$. Indeed, assume that $\rho(\hat{\theta}, \theta_v) < \delta$; then for any $v' \neq v$, we have

$$\rho(\hat{\theta}, \theta_{v'}) \geq \rho(\theta_v, \theta_{v'}) - \rho(\hat{\theta}, \theta_v) > 2\delta - \delta = \delta.$$

The test must thus return v as claimed. Equivalently, for $v \in \mathcal{V}$, the inequality $\Psi(\hat{\theta}) \neq v$ implies $\rho(\hat{\theta}, \theta_v) \geq \delta$. (See Figure 9.1.) By averaging over \mathcal{V} , we find that

$$\sup_P \mathbb{P}(\rho(\hat{\theta}, \theta(P)) \geq \delta) \geq \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \mathbb{P}(\rho(\hat{\theta}, \theta(P_v)) \geq \delta \mid V = v) \geq \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \mathbb{P}(\Psi(\hat{\theta}) \neq v \mid V = v).$$

Taking an infimum over all tests $\Psi : \mathcal{X}^n \rightarrow V$ gives inequality (9.2.1). \square

The remaining challenge is to lower bound the probability of error in the underlying multi-way hypothesis testing problem, which we do by choosing the separation δ to trade off between the loss $\Phi(\delta)$ (large δ increases the loss) and the probability of error (small δ , and hence separation, makes the hypothesis test harder). Usually, one attempts to choose the largest separation δ that guarantees a constant probability of error. There are a variety of techniques for this, and we present three: Le Cam's method, Fano's method, and Assouad's method, including extensions of the latter two to enhance their applicability. Before continuing, however, we review some inequalities between divergence measures defined on probabilities, which will be essential for our development, and concepts related to packing sets (metric entropy, covering numbers, and packing).

9.2.2 Inequalities between divergences and product distributions

We now present a few inequalities, and their consequences when applied to product distributions, that will be quite useful for proving our lower bounds. The three divergences we relate are the total variation distance, Kullback-Leibler divergence, and Hellinger distance, all of which are instances

of f -divergences (recall Section 2.2.3). We first recall the definitions of the three when applied to distributions P, Q on a set \mathcal{X} , which we assume have densities p, q with respect to a base measure μ . Then we recall the total variation distance (2.2.6) is

$$\|P - Q\|_{\text{TV}} := \sup_{A \subset \mathcal{X}} |P(A) - Q(A)| = \frac{1}{2} \int |p(x) - q(x)| d\mu(x),$$

which is the f -divergence $D_f(P\|Q)$ generated by $f(t) = \frac{1}{2}|t - 1|$. The Hellinger distance (2.2.7) is

$$d_{\text{hel}}(P, Q)^2 := \int (\sqrt{p(x)} - \sqrt{q(x)})^2 d\mu(x),$$

which is the f -divergence $D_f(P\|Q)$ generated by $f(t) = (\sqrt{t} - 1)^2$. We also recall the Kullback-Leibler (KL) divergence

$$D_{\text{kl}}(P\|Q) := \int p(x) \log \frac{p(x)}{q(x)} d\mu(x), \quad (9.2.2)$$

which is the f -divergence $D_f(P\|Q)$ generated by $f(t) = t \log t$. As noted in Section 2.2.3, Proposition 2.2.8, these divergences have the following relationships.

Proposition (Proposition 2.2.8, restated). *The total variation distance satisfies the following relationships:*

(a) *For the Hellinger distance,*

$$\frac{1}{2} d_{\text{hel}}(P, Q)^2 \leq \|P - Q\|_{\text{TV}} \leq d_{\text{hel}}(P, Q) \sqrt{1 - d_{\text{hel}}(P, Q)^2/4}.$$

(b) *Pinsker's inequality: for any distributions P, Q ,*

$$\|P - Q\|_{\text{TV}}^2 \leq \frac{1}{2} D_{\text{kl}}(P\|Q).$$

We now show how Proposition 2.2.8 is useful, because KL-divergence and Hellinger distance both are easier to manipulate on product distributions than is total variation. Specifically, consider the product distributions $P = P_1 \times \cdots \times P_n$ and $Q = Q_1 \times \cdots \times Q_n$. Then the KL-divergence satisfies the decoupling equality

$$D_{\text{kl}}(P\|Q) = \sum_{i=1}^n D_{\text{kl}}(P_i\|Q_i), \quad (9.2.3)$$

while because $d_{\text{hel}}^2(P, Q) = 1 - \int \sqrt{dP dQ}$, the Hellinger distance satisfies

$$\begin{aligned} d_{\text{hel}}^2(P, Q) &= 1 - \int \sqrt{p_1(x_1) \cdots p_n(x_n) q_1(x_1) \cdots q_n(x_n)} d\mu(x_1^n) \\ &= 1 - \prod_{i=1}^n \int \sqrt{p_i(x_i) q_i(x_i)} d\mu(x_i) \\ &= 1 - \prod_{i=1}^n (1 - d_{\text{hel}}^2(P_i, Q_i)). \end{aligned} \quad (9.2.4)$$

In particular, we see that for product distributions P^n and Q^n , Proposition 2.2.8 implies that

$$\|P^n - Q^n\|_{\text{TV}}^2 \leq \frac{1}{2} D_{\text{kl}}(P^n \| Q^n) = \frac{n}{2} D_{\text{kl}}(P \| Q)$$

and

$$\|P^n - Q^n\|_{\text{TV}} \leq d_{\text{hel}}(P^n, Q^n) \leq \sqrt{2 - 2(1 - d_{\text{hel}}(P, Q))^2}^n.$$

As a consequence, if we can guarantee that $D_{\text{kl}}(P \| Q) \leq 1/n$ or $d_{\text{hel}}(P, Q) \leq 1/\sqrt{n}$, then we guarantee the strict inequality $\|P^n - Q^n\|_{\text{TV}} \leq 1 - c$ for a fixed constant $c > 0$, for any n . We will see how this type of guarantee can be used to prove minimax lower bounds in the following sections.

9.2.3 Metric entropy and packing numbers

The second part of proving our lower bounds involves the construction of the packing set in Section 9.2.1. The size of the space Θ of parameters associated with our estimation problem—and consequently, how many parameters we can pack into it—is strongly coupled with the difficulty of estimation. The tools we develop in Section 5.1.2 on metric entropies and covering and packing numbers therefore become central.

Probably the most central construction relies on volume bounds on packing and covering numbers, which we recall from Lemma 5.1.10: the covering and packing numbers of a norm ball \mathbb{B} in its own norm $\|\cdot\|$ scale exponentially in the dimension. In particular, for any $\delta < 1$, there is a packing \mathcal{V} of \mathbb{B} such that $\|v - v'\| \geq \delta$ for all distinct $v, v' \in \mathcal{V}$ and $|\mathcal{V}| \geq (1/\delta)^d$, because we know $M(\delta, \mathbb{B}, \|\cdot\|) \geq N(\delta, \mathbb{B}, \|\cdot\|)$ as in Lemma 5.1.8. We thus state the following corollary for later use, which states that we can construct exponentially large packings of arbitrary norm-balls (in finite dimensions) where the points have constant distance from one another.

Corollary 9.2.2. *Let $\mathbb{B}^d = \{v \in \mathbb{R}^d \mid \|v\| \leq 1\}$ be the unit ball for the norm $\|\cdot\|$. Then there exists $\mathcal{V} \subset \mathbb{B}^d$ with $|\mathcal{V}| \geq 2^d$ and $\|v - v'\| \geq \frac{1}{2}$ for each $v \neq v' \in \mathcal{V}$.*

Another common packing arises from coding theory, where the technique is to construct well-separated code-books ($\{0, 1\}$ -valued bit strings associated to individual symbols to be communicated) for communication. In showing our lower bounds, we show that even if a code-book is well-separated, it may still be hard to estimate. With that, we now demonstrate that there exist (exponentially) large packings of the d -dimensional hypercube of points that are $O(d)$ -separated in the Hamming metric.

Lemma 9.2.3 (Gilbert-Varshamov bound). *Let $d \geq 1$. There is a subset \mathcal{V} of the d -dimensional hypercube $\mathcal{H}_d = \{-1, 1\}^d$ of size $|\mathcal{V}| \geq \exp(d/8)$ such that the ℓ_1 -distance*

$$\|v - v'\|_1 = 2 \sum_{j=1}^d \mathbf{1}\{v_j \neq v'_j\} \geq \frac{d}{2}$$

for all $v \neq v'$ with $v, v' \in \mathcal{V}$.

Proof We use the proof of Guntuboyina [106]. Consider a maximal subset \mathcal{V} of $\mathcal{H}_d = \{-1, 1\}^d$ satisfying

$$\|v - v'\|_1 \geq d/2 \quad \text{for all distinct } v, v' \in \mathcal{V}. \quad (9.2.5)$$

That is, the addition of any vector $w \in \mathcal{H}_d, w \notin \mathcal{V}$ to \mathcal{V} will break the constraint (9.2.5). This means that if we construct the closed balls $B(v, d/2) := \{w \in \mathcal{H}_d : \|v - w\|_1 \leq d/2\}$, we must have

$$\bigcup_{v \in \mathcal{V}} B(v, d/2) = \mathcal{H}_d \quad \text{so} \quad |\mathcal{V}| |B(0, d/2)| = \sum_{v \in \mathcal{V}} |B(v, d/2)| \geq 2^d. \quad (9.2.6)$$

We now upper bound the cardinality of $B(v, d/2)$ using the probabilistic method, which will imply the desired result. Let $S_i, i = 1, \dots, d$, be i.i.d. Bernoulli $\{0, 1\}$ -valued random variables. Then by their uniformity, for any $v \in \mathcal{H}_d$,

$$\begin{aligned} 2^{-d} |B(v, d/2)| &= \mathbb{P}(S_1 + S_2 + \dots + S_d \leq d/4) = \mathbb{P}(S_1 + S_2 + \dots + S_d \geq 3d/4) \\ &\leq \mathbb{E}[\exp(\lambda S_1 + \dots + \lambda S_d)] \exp(-3\lambda d/4) \end{aligned}$$

for any $\lambda > 0$, by Markov's inequality (or the Chernoff bound). Since $\mathbb{E}[\exp(\lambda S_1)] = \frac{1}{2}(1 + e^\lambda)$, we obtain

$$2^{-d} |B(v, d/2)| \leq \inf_{\lambda \geq 0} \left\{ 2^{-d} (1 + e^\lambda)^d \exp(-3\lambda d/4) \right\}$$

Choosing $\lambda = \log 3$, we have

$$|B(v, d/2)| \leq 4^d \exp(-(3/4)d \log 3) = 3^{-3d/4} 4^d.$$

Recalling inequality (9.2.6), we have

$$|\mathcal{V}| 3^{-3d/4} 4^d \geq |\mathcal{V}| |B(v, d/2)| \geq 2^d, \quad \text{or} \quad |\mathcal{V}| \geq \frac{3^{3d/4}}{2^d} = \exp\left(d \left[\frac{3}{4} \log 3 - \log 2 \right]\right) \geq \exp(d/8),$$

as claimed. \square

9.3 Le Cam's method

Le Cam's method, in its simplest form, provides lower bounds on the error in simple binary hypothesis testing problems. In this section, we explore this connection, showing the connection between hypothesis testing and total variation distance, and we then show how this can yield lower bounds on minimax error (or the optimal Bayes' risk) for simple—often one-dimensional—estimation problems.

In the first homework, we considered several representations of the total variation distance, including a question showing its relation to optimal testing. We begin again with this strand of thought, recalling the general testing problem discussed in Section 9.2.1. Suppose that we have a Bayesian hypothesis testing problem where V is chosen with equal probability to be 1 or 2, and given $V = v$, the sample X is drawn from the distribution P_v . Denoting by \mathbb{P} the joint distribution of V and X , we have for any test $\Psi : \mathcal{X} \rightarrow \{1, 2\}$ that the probability of error is

$$\mathbb{P}(\Psi(X) \neq V) = \frac{1}{2} P_1(\Psi(X) \neq 1) + \frac{1}{2} P_2(\Psi(X) \neq 2).$$

Recalling Section 9.2.1, we note that Proposition 2.3.1 gives an exact representation of the testing error using total variation distance. In particular, we have

Proposition (Proposition 2.3.1, restated). *For any distributions P_1 and P_2 on \mathcal{X} , we have*

$$\inf_{\Psi} \{P_1(\Psi(X) \neq 1) + P_2(\Psi(X) \neq 2)\} = 1 - \|P_1 - P_2\|_{\text{TV}}, \quad (9.3.1)$$

where the infimum is taken over all tests $\Psi : \mathcal{X} \rightarrow \{1, 2\}$.

Returning to the setting in which we receive n i.i.d. observations $X_i \sim P$, when $V = 1$ with probability $\frac{1}{2}$ and 2 with probability $\frac{1}{2}$, we have

$$\inf_{\Psi} \mathbb{P}(\Psi(X_1, \dots, X_n) \neq V) = \frac{1}{2} - \frac{1}{2} \|P_1^n - P_2^n\|_{\text{TV}}. \quad (9.3.2)$$

The representations (9.3.1) and (9.3.2), in conjunction with our reduction of estimation to testing in Proposition 9.2.1, imply the following lower bound on minimax risk. For any family \mathcal{P} of distributions for which there exists a pair $P_1, P_2 \in \mathcal{P}$ satisfying $\rho(\theta(P_1), \theta(P_2)) \geq 2\delta$, then the minimax risk after n observations has lower bound

$$\mathfrak{M}_n(\theta(\mathcal{P}), \Phi \circ \rho) \geq \Phi(\delta) \left[\frac{1}{2} - \frac{1}{2} \|P_1^n - P_2^n\|_{\text{TV}} \right]. \quad (9.3.3)$$

The lower bound (9.3.3) suggests the following strategy: we find distributions P_1 and P_2 , which we choose as a function of δ , that guarantee $\|P_1^n - P_2^n\|_{\text{TV}} \leq \frac{1}{2}$. In this case, so long as $\rho(\theta(P_1), \theta(P_2)) \geq 2\delta$, we have the lower bound

$$\mathfrak{M}_n(\theta(\mathcal{P}), \Phi \circ \rho) \geq \Phi(\delta) \left[\frac{1}{2} - \frac{1}{2} \cdot \frac{1}{4} \right] = \frac{1}{4} \Phi(\delta).$$

We now give an example illustrating this idea.

Example 9.3.1 (Bernoulli mean estimation): Consider the problem of estimating the mean $\theta \in [-1, 1]$ of a $\{\pm 1\}$ -valued Bernoulli distribution under the squared error loss $(\theta - \hat{\theta})^2$, where $X_i \in \{-1, 1\}$. In this case, by fixing some $\delta > 0$, we set $\mathcal{V} = \{-1, 1\}$, and we define P_v so that

$$P_v(X = 1) = \frac{1 + v\delta}{2} \quad \text{and} \quad P_v(X = -1) = \frac{1 - v\delta}{2},$$

whence we see that the mean $\theta(P_v) = \delta v$. Using the metric $\rho(\theta, \theta') = |\theta - \theta'|$ and loss $\Phi(\delta) = \delta^2$, we have separation 2δ of $\theta(P_{-1})$ and $\theta(P_1)$. Thus, via Le Cam's method (9.3.3), we have that

$$\mathfrak{M}_n(\text{Bernoulli}([-1, 1]), (\cdot)^2) \geq \frac{1}{2} \delta^2 (1 - \|P_{-1}^n - P_1^n\|_{\text{TV}}).$$

We would thus like to upper bound $\|P_{-1}^n - P_1^n\|_{\text{TV}}$ as a function of the separation δ and sample size n ; here we use Pinsker's inequality (Proposition 2.2.8(a)) and the tensorization identity (9.2.3) that makes KL-divergence so useful. Indeed, we have

$$\|P_{-1}^n - P_1^n\|_{\text{TV}}^2 \leq \frac{1}{2} D_{\text{kl}}(P_{-1}^n \| P_1^n) = \frac{n}{2} D_{\text{kl}}(P_{-1} \| P_1) = \frac{n}{2} \delta \log \frac{1 + \delta}{1 - \delta}.$$

Noting that $\delta \log \frac{1 + \delta}{1 - \delta} \leq 3\delta^2$ for $\delta \in [0, 1/2]$, we obtain that $\|P_{-1}^n - P_1^n\|_{\text{TV}} \leq \delta \sqrt{3n/2}$ for $\delta \leq 1/2$. In particular, we can guarantee a high probability of error in the associated hypothesis testing problem (recall inequality (9.3.2)) by taking $\delta = 1/\sqrt{6n}$; this guarantees $\|P_{-1}^n - P_1^n\|_{\text{TV}} \leq \frac{1}{2}$. We thus have the minimax lower bound

$$\mathfrak{M}_n(\text{Bernoulli}([-1, 1]), (\cdot)^2) \geq \frac{1}{2} \delta^2 \left(1 - \frac{1}{2} \right) = \frac{1}{24n}.$$

While the factor $1/24$ is smaller than necessary, this bound is optimal to within constant factors; the sample mean $(1/n) \sum_{i=1}^n X_i$ achieves mean-squared error $(1 - \theta^2)/n$.

As an alternative proof, we may use the Hellinger distance and its associated decoupling identity (9.2.4). We sketch the idea, ignoring lower order terms when convenient. In this case, Proposition 2.2.7 implies

$$\|P_1^n - P_2^n\|_{\text{TV}} \leq \sqrt{2} d_{\text{hel}}(P_1^n, P_2^n) = \sqrt{2 - 2(1 - d_{\text{hel}}(P_1, P_2)^2)^n}.$$

Noting that

$$d_{\text{hel}}(P_1, P_2)^2 = \left(\sqrt{\frac{1+\delta}{2}} - \sqrt{\frac{1-\delta}{2}} \right)^2 = 1 - 2\sqrt{\frac{1-\delta^2}{4}} = 1 - \sqrt{1-\delta^2} \approx \frac{1}{2}\delta^2,$$

and noting that $(1 - \delta^2) \approx e^{-\delta^2}$, we have (up to lower order terms in δ) that $\|P_1^n - P_2^n\|_{\text{TV}} \leq \sqrt{2 - 2\exp(-\delta^2 n/2)}$. Choosing $\delta^2 = 1/(4n)$, we have $\sqrt{2 - 2\exp(-\delta^2 n/2)} \leq 1/2$, thus giving the lower bound

$$\mathfrak{M}_n(\text{Bernoulli}([-1, 1]), (\cdot)^2) \text{ “} \geq \text{” } \frac{1}{2}\delta^2 \left(1 - \frac{1}{2}\right) = \frac{1}{16n},$$

where the quotations indicate we have been fast and loose in the derivation. \diamond

This example shows the “usual” rate of convergence in parametric estimation problems, that is, that we can estimate a parameter θ at a rate (in squared error) scaling as $1/n$. The mean estimator above is, in some sense, the prototypical example of such regular problems. In some “irregular” scenarios—including estimating the support of a uniform random variable, which we study in the homework—faster rates are possible.

We also note in passing that there are substantially more complex versions of Le Cam’s method that can yield sharp results for a wider variety of problems, including some in nonparametric estimation [134, 194]. For our purposes, the simpler two-point perspective provided in this section will be sufficient.

JCD Comment: Talk about Euclidean structure with KL space and information geometry a bit here to suggest the KL approach later.

9.4 Fano’s method

Fano’s method, originally proposed by Has’minskii [109] for providing lower bounds in nonparametric estimation problems, gives a somewhat more general technique than Le Cam’s method, and it applies when the packing set \mathcal{V} has cardinality larger than two. The method has played a central role in minimax theory, beginning with the pioneering work of Has’minskii and Ibragimov [109, 115]. More recent work following this initial push continues to the present day (e.g. [31, 194, 192, 32, 160, 106, 47]).

9.4.1 The classical (local) Fano method

We begin by stating Fano’s inequality, which provides a lower bound on the error in a multi-way hypothesis testing problem. Let V be a random variable taking values in a finite set \mathcal{V} with cardinality $|\mathcal{V}| \geq 2$. If we let the function $h_2(p) = -p \log p - (1 - p) \log(1 - p)$ denote the entropy of the Bernoulli random variable with parameter p , Fano’s inequality (Proposition 2.3.3 from Chapter 2) takes the following form:

Proposition 9.4.1 (Fano inequality). *For any Markov chain $V \rightarrow X \rightarrow \hat{V}$, we have*

$$h_2(\mathbb{P}(\hat{V} \neq V)) + \mathbb{P}(\hat{V} \neq V) \log(|\mathcal{V}| - 1) \geq H(V | \hat{V}). \quad (9.4.1)$$

Restating the results in Chapter 2, we also have the following convenient rewriting of Fano's inequality when V is uniform in \mathcal{V} (recall Corollary 2.3.4).

Corollary 9.4.2. *Assume that V is uniform on \mathcal{V} . For any Markov chain $V \rightarrow X \rightarrow \hat{V}$,*

$$\mathbb{P}(\hat{V} \neq V) \geq 1 - \frac{I(V; X) + \log 2}{\log(|\mathcal{V}|)}. \quad (9.4.2)$$

In particular, Corollary 9.4.2 shows that we have

$$\inf_{\Psi} \mathbb{P}(\Psi(X) \neq V) \geq 1 - \frac{I(V; X) + \log 2}{\log |\mathcal{V}|},$$

where the infimum is taken over all testing procedures Ψ . By combining Corollary 9.4.2 with the reduction from estimation to testing in Proposition 9.2.1, we obtain the following result.

Proposition 9.4.3. *Let $\{\theta(P_v)\}_{v \in \mathcal{V}}$ be a 2δ -packing in the ρ -semimetric. Assume that V is uniform on the set \mathcal{V} , and conditional on $V = v$, we draw a sample $X \sim P_v$. Then the minimax risk has lower bound*

$$\mathfrak{M}(\theta(\mathcal{P}); \Phi \circ \rho) \geq \Phi(\delta) \left(1 - \frac{I(V; X) + \log 2}{\log |\mathcal{V}|} \right).$$

To gain some intuition for Proposition 9.4.3, we think of the lower bound as a function of the separation $\delta > 0$. Roughly, as $\delta \downarrow 0$, the separation condition between the distributions P_v is relaxed and we expect the distributions P_v to be closer to one another. In this case—as will be made more explicitly presently—the hypothesis testing problem of distinguishing the P_v becomes more challenging, and the information $I(V; X)$ shrinks. Thus, what we roughly attempt to do is to choose our packing $\theta(P_v)$ as a function of δ , and find the largest $\delta > 0$ making the mutual information small enough that

$$\frac{I(V; X) + \log 2}{\log |\mathcal{V}|} \leq \frac{1}{2}. \quad (9.4.3)$$

In this case, the minimax lower bound is at least $\Phi(\delta)/2$. We now explore techniques for achieving such results.

Mutual information and KL-divergence

Many techniques for upper bounding mutual information rely on its representation as the KL-divergence between multiple distributions. Indeed, given random variables V and X as in the preceding sections, if we let $P_{V,X}$ denote their joint distribution and P_V and P_X their marginals, then

$$I(V; X) = D_{\text{kl}}(P_{X,V} \| P_X \times P_V),$$

where $P_X \times P_V$ denotes the distribution of (X, V) when the random variables are independent. By manipulating this definition, we can rewrite it into a form more convenient for our purposes.

Indeed, focusing on our setting of testing, let us assume that V is drawn from a prior distribution π (this may be a discrete or arbitrary distribution, though for simplicity we focus on the case when

π is discrete). Let P_v denote the distribution of X conditional on $V = v$, as in Proposition 9.4.3. Then marginally, we know that X is drawn from the mixture distribution

$$\bar{P} := \sum_v \pi(v) P_v.$$

With this definition of the mixture distribution, via algebraic manipulations, we have

$$I(V; X) = \sum_v \pi(v) D_{\text{kl}}(P_v \| \bar{P}), \quad (9.4.4)$$

a representation that plays an important role in our subsequent derivations. To see equality (9.4.4), let μ be a base measure over \mathcal{X} (assume w.l.o.g. that X has density $p(\cdot | v) = p_v(\cdot)$ conditional on $V = v$), and note that

$$I(V; X) = \sum_v \int_{\mathcal{X}} p(x | v) \pi(v) \log \frac{p(x | v)}{\sum_{v'} p(x | v') \pi(v')} d\mu(x) = \sum_v \pi(v) \int_{\mathcal{X}} p(x | v) \log \frac{p(x | v)}{\bar{p}(x)} d\mu(x).$$

Representation (9.4.4) makes it clear that if the distributions of the sample X conditional on V are all similar, then there is little information content. Returning to the discussion after Proposition 9.4.3, we have in this uniform setting that

$$\bar{P} = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P_v \quad \text{and} \quad I(V; X) = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} D_{\text{kl}}(P_v \| \bar{P}).$$

The mutual information is small if the typical conditional distribution P_v is difficult to distinguish—has small KL-divergence—from \bar{P} .

The local Fano method

The local Fano method is based on a weakening of the mixture representation of mutual information (9.4.4), then giving a uniform upper bound on divergences between all pairs of the conditional distributions P_v and $P_{v'}$. (This method is known in the statistics literature as the “generalized Fano method,” a poor name, as it is based on a weak upper bound on mutual information.) In particular (focusing on the case when V is uniform), the convexity of $-\log$ implies that

$$I(V; X) = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} D_{\text{kl}}(P_v \| \bar{P}) \leq \frac{1}{|\mathcal{V}|^2} \sum_{v, v'} D_{\text{kl}}(P_v \| P_{v'}). \quad (9.4.5)$$

In the local Fano method approach, we construct a *local packing*. This local packing approach is based on constructing a family of distributions P_v for $v \in \mathcal{V}$ defining a 2δ -packing (recall Section 9.2.1), meaning that $\rho(\theta(P_v), \theta(P_{v'})) \geq 2\delta$ for all $v \neq v'$, but which additionally satisfy the uniform upper bound

$$D_{\text{kl}}(P_v \| P_{v'}) \leq \kappa^2 \delta^2 \quad \text{for all } v, v' \in \mathcal{V}, \quad (9.4.6)$$

where $\kappa > 0$ is a fixed problem-dependent constant. If we have the inequality (9.4.6), then so long as we can find a *local* packing \mathcal{V} such that

$$\log |\mathcal{V}| \geq 2(\kappa^2 \delta^2 + \log 2),$$

we are guaranteed the testing error condition (9.4.3), and hence the minimax lower bound

$$\mathfrak{M}(\theta(\mathcal{P}), \Phi \circ \rho) \geq \frac{1}{2} \Phi(\delta).$$

The difficulty in this approach is constructing the packing set \mathcal{V} that allows δ to be chosen to obtain sharp lower bounds, and we often require careful choices of the packing sets \mathcal{V} . (We will see how to reduce such difficulties in subsequent sections.)

Constructing local packings As mentioned above, the main difficulty in using Fano's method is in the construction of so-called “local” packings. In these problems, the idea is to construct a packing \mathcal{V} of a fixed set (in a vector space, say \mathbb{R}^d) with constant radius and constant distance. Then we scale elements of the packing by $\delta > 0$, which leaves the cardinality $|\mathcal{V}|$ identical, but allows us to scale δ in the separation in the packing and the uniform divergence bound (9.4.6). In particular, Lemmas 9.2.3 and 5.1.10 show that we can construct exponentially large packings of certain sets with balls of a fixed radius.

We now illustrate these techniques via two examples.

Example 9.4.4 (Normal mean estimation): Consider the d -dimensional normal location family $\mathcal{N}_d = \{\mathbf{N}(\theta, \sigma^2 I_{d \times d}) \mid \theta \in \mathbb{R}^d\}$; we wish to estimate the mean $\theta = \theta(P)$ of a given distribution $P \in \mathcal{N}_d$ in mean-squared error, that is, with loss $\|\hat{\theta} - \theta\|_2^2$. Let \mathcal{V} be a $1/2$ -packing of the unit ℓ_2 -ball with cardinality at least 2^d , as guaranteed by Lemma 5.1.10. (We assume for simplicity that $d \geq 2$.)

Now we construct our local packing. Fix $\delta > 0$, and for each $v \in \mathcal{V}$, set $\theta_v = \delta v \in \mathbb{R}^d$. Then we have

$$\|\theta_v - \theta_{v'}\|_2 = \delta \|v - v'\|_2 \geq \frac{\delta}{2}$$

for each distinct pair $v, v' \in \mathcal{V}$, and moreover, we note that $\|\theta_v - \theta_{v'}\|_2 \leq \delta$ for such pairs as well. By applying the Fano minimax bound of Proposition 9.4.3, we see that (given n normal observations $X_i \stackrel{\text{iid}}{\sim} P$)

$$\mathfrak{M}_n(\theta(\mathcal{N}_d), \|\cdot\|_2^2) \geq \left(\frac{1}{2} \cdot \frac{\delta}{2}\right)^2 \left(1 - \frac{I(V; X_1^n) + \log 2}{\log |\mathcal{V}|}\right) = \frac{\delta^2}{16} \left(1 - \frac{I(V; X_1^n) + \log 2}{d \log 2}\right).$$

Now note that for any pair v, v' , if P_v is the normal distribution $\mathbf{N}(\theta_v, \sigma^2 I_{d \times d})$ we have

$$D_{\text{kl}}(P_v^n \| P_{v'}^n) = n \cdot D_{\text{kl}}(\mathbf{N}(\delta v, \sigma^2 I_{d \times d}) \| \mathbf{N}(\delta v', \sigma^2 I_{d \times d})) = n \cdot \frac{\delta^2}{2\sigma^2} \|v - v'\|_2^2,$$

as the KL-divergence between two normal distributions with identical covariance is

$$D_{\text{kl}}(\mathbf{N}(\theta_1, \Sigma) \| \mathbf{N}(\theta_2, \Sigma)) = \frac{1}{2} (\theta_1 - \theta_2)^\top \Sigma^{-1} (\theta_1 - \theta_2)$$

as in Example 2.1.7. As $\|v - v'\|_2 \leq 1$, we have the KL-divergence bound (9.4.6) with $\kappa^2 = n/2\sigma^2$.

Combining our derivations, we have the minimax lower bound

$$\mathfrak{M}_n(\theta(\mathcal{N}_d), \|\cdot\|_2^2) \geq \frac{\delta^2}{16} \left(1 - \frac{n\delta^2/2\sigma^2 + \log 2}{d \log 2}\right). \quad (9.4.7)$$

Then by taking $\delta^2 = d\sigma^2 \log 2 / (2n)$, we see that

$$1 - \frac{n\delta^2/2\sigma^2 + \log 2}{d \log 2} = 1 - \frac{1}{d} - \frac{1}{4} \geq \frac{1}{4}$$

by assumption that $d \geq 2$, and inequality (9.4.7) implies the minimax lower bound

$$\mathfrak{M}_n(\theta(\mathcal{N}_d), \|\cdot\|_2^2) \geq \frac{d\sigma^2 \log 2}{32n} \cdot \frac{1}{4} \geq \frac{1}{185} \cdot \frac{d\sigma^2}{n}.$$

While the constant $1/185$ is not sharp, we do obtain the right scaling in d , n , and the variance σ^2 ; the sample mean attains the same risk. \diamond

Example 9.4.5 (Linear regression): In this example, we show how local packings can give (up to some constant factors) sharp minimax rates for standard linear regression problems. In particular, for fixed matrix $X \in \mathbb{R}^{n \times d}$, we observe

$$Y = X\theta + \varepsilon,$$

where $\varepsilon \in \mathbb{R}^n$ consists of independent random variables ε_i with variance bounded by $\text{Var}(\varepsilon_i) \leq \sigma^2$, and $\theta \in \mathbb{R}^d$ is allowed to vary over \mathbb{R}^d . For the purposes of our lower bound, we may assume that $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_{n \times n})$. Let \mathcal{P} denote the family of such normally distributed linear regression problems, and assume for simplicity that $d \geq 32$.

In this case, we use the Gilbert-Varshamov bound (Lemma 9.2.3) to construct a local packing and attain minimax rates. Indeed, let \mathcal{V} be a packing of $\{-1, 1\}^d$ such that $\|v - v'\|_1 \geq d/2$ for distinct elements of \mathcal{V} , and let $|\mathcal{V}| \geq \exp(d/8)$ as guaranteed by the Gilbert-Varshamov bound. For fixed $\delta > 0$, if we set $\theta_v = \delta v$, then we have the packing guarantee for distinct elements v, v' that

$$\|\theta_v - \theta_{v'}\|_2^2 = \delta^2 \sum_{j=1}^d (v_j - v'_j)^2 = 4\delta^2 \|v - v'\|_1 \geq 2d\delta^2.$$

Moreover, we have the upper bound

$$\begin{aligned} D_{\text{kl}}(\mathcal{N}(X\theta_v, \sigma^2 I_{n \times n}) \| \mathcal{N}(X\theta_{v'}, \sigma^2 I_{n \times n})) &= \frac{1}{2\sigma^2} \|X(\theta_v - \theta_{v'})\|_2^2 \\ &\leq \frac{\delta^2}{2\sigma^2} \gamma_{\max}^2(X) \|v - v'\|_2^2 \leq \frac{2d}{\sigma^2} \gamma_{\max}^2(X) \delta^2, \end{aligned}$$

where $\gamma_{\max}(X)$ denotes the maximum singular value of X . Consequently, the bound (9.4.6) holds with $\kappa^2 \leq 2d\gamma_{\max}^2(X)/\sigma^2$, and we have the minimax lower bound

$$\mathfrak{M}(\theta(\mathcal{P}), \|\cdot\|_2^2) \geq \frac{d\delta^2}{2} \left(1 - \frac{I(V; Y) + \log 2}{\log |\mathcal{V}|} \right) \geq \frac{d\delta^2}{2} \left(1 - \frac{\frac{2d\gamma_{\max}^2(X)}{\sigma^2} \delta^2 + \log 2}{d/8} \right).$$

Now, if we choose

$$\delta^2 = \frac{\sigma^2}{64\gamma_{\max}^2(X)}, \quad \text{then} \quad 1 - \frac{8 \log 2}{d} - \frac{16d\gamma_{\max}^2(X)\delta^2}{d} \geq 1 - \frac{1}{4} - \frac{1}{4} = \frac{1}{2},$$

by assumption that $d \geq 32$. In particular, we obtain the lower bound

$$\mathfrak{M}(\theta(\mathcal{P}), \|\cdot\|_2^2) \geq \frac{1}{256} \frac{\sigma^2 d}{\gamma_{\max}^2(X)} = \frac{1}{256} \frac{\sigma^2 d}{n} \frac{1}{\gamma_{\max}^2(\frac{1}{\sqrt{n}}X)},$$

for a convergence rate (roughly) of $\sigma^2 d/n$ after rescaling the singular values of X by $1/\sqrt{n}$.

This bound is sharp in terms of the dimension, dependence on n , and the variance σ^2 , but it does not fully capture the dependence on X , as it depends only on the maximum singular value. Indeed, in this case, an exact calculation (cf. [136]) shows that the minimax value of the problem is exactly $\sigma^2 \text{tr}((X^\top X)^{-1})$. Letting $\lambda_j(A)$ be the j th eigenvalue of a matrix A , we have

$$\begin{aligned} \sigma^2 \text{tr}((X^\top X)^{-1}) &= \frac{\sigma^2}{n} \text{tr}((n^{-1} X^\top X)^{-1}) = \frac{\sigma^2}{n} \sum_{j=1}^d \frac{1}{\lambda_j(\frac{1}{n} X^\top X)} \\ &\geq \frac{\sigma^2 d}{n} \min_j \frac{1}{\lambda_j(\frac{1}{n} X^\top X)} = \frac{\sigma^2 d}{n} \frac{1}{\gamma_{\max}^2(\frac{1}{\sqrt{n}}X)}. \end{aligned}$$

Thus, the local Fano method captures most—but not all—of the difficulty of the problem. \diamond

9.4.2 A distance-based Fano method

While the testing lower bound (9.4.2) is sufficient for proving lower bounds for many estimation problems, for the sharpest results it sometimes requires a somewhat delicate construction of a well-separated packing (e.g. [47, 75]). To that end, we also provide extensions of inequalities (9.4.1) and (9.4.2) that more directly yield bounds on estimation error, allowing more direct and simpler proofs of a variety of minimax lower bounds (see also reference [72]).

More specifically, suppose that the distance function $\rho_{\mathcal{V}}$ is defined on \mathcal{V} , and we are interested in bounding the estimation error $\rho_{\mathcal{V}}(\hat{V}, V)$. We begin by providing analogues of the lower bounds (9.4.1) and (9.4.2) that replace the testing error with the tail probability $\mathbb{P}(\rho_{\mathcal{V}}(\hat{V}, V) > t)$. By Markov's inequality, such control directly yields bounds on the expectation $\mathbb{E}[\rho_{\mathcal{V}}(\hat{V}, V)]$. As we show in the sequel and in chapters to come, these distance-based Fano inequalities allow more direct proofs of a variety of minimax bounds without the need for careful construction of packing sets or metric entropy calculations as in other arguments.

We begin with the distance-based analogue of the usual discrete Fano inequality in Proposition 9.4.1. Let V be a random variable supported on a finite set \mathcal{V} with cardinality $|\mathcal{V}| \geq 2$, and let $\rho : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$ be a function defined on $\mathcal{V} \times \mathcal{V}$. In the usual setting, the function ρ is a metric on the space \mathcal{V} , but our theory applies to general functions. For a given scalar $t \geq 0$, the maximum and minimum *neighborhood sizes at radius t* are given by

$$N_t^{\max} := \max_{v \in \mathcal{V}} \{\text{card}\{v' \in \mathcal{V} \mid \rho(v, v') \leq t\}\} \quad \text{and} \quad N_t^{\min} := \min_{v \in \mathcal{V}} \{\text{card}\{v' \in \mathcal{V} \mid \rho(v, v') \leq t\}\}. \quad (9.4.8)$$

Defining the error probability $P_t = \mathbb{P}(\rho_{\mathcal{V}}(\hat{V}, V) > t)$, we then have the following generalization of Fano's inequality:

Proposition 9.4.6. *For any Markov chain $V \rightarrow X \rightarrow \hat{V}$, we have*

$$h_2(P_t) + P_t \log \frac{|\mathcal{V}| - N_t^{\min}}{N_t^{\max}} + \log N_t^{\max} \geq H(V \mid \hat{V}). \quad (9.4.9)$$

Before proving the proposition, which we do in Section 9.6.1, it is informative to note that it reduces to the standard form of Fano's inequality (9.4.1) in a special case. Suppose that we take $\rho_{\mathcal{V}}$ to be the 0-1 metric, meaning that $\rho_{\mathcal{V}}(v, v') = 0$ if $v = v'$ and 1 otherwise. Setting $t = 0$ in Proposition 9.4.6, we have $P_0 = \mathbb{P}[\hat{V} \neq V]$ and $N_0^{\min} = N_0^{\max} = 1$, whence inequality (9.4.9) reduces to inequality (9.4.1). Other weakenings allow somewhat clearer statements (see Section 9.6.2 for a proof):

Corollary 9.4.7. *If V is uniform on \mathcal{V} and $(|\mathcal{V}| - N_t^{\min}) > N_t^{\max}$, then*

$$\mathbb{P}(\rho_{\mathcal{V}}(\hat{V}, V) > t) \geq 1 - \frac{I(V; X) + \log 2}{\log \frac{|\mathcal{V}|}{N_t^{\max}}}. \quad (9.4.10)$$

Inequality (9.4.10) is the natural analogue of the classical mutual-information based form of Fano's inequality (9.4.2), and it provides a qualitatively similar bound. The main difference is that the usual cardinality $|\mathcal{V}|$ is replaced by the ratio $|\mathcal{V}|/N_t^{\max}$. This quantity serves as a rough measure of the number of possible “regions” in the space \mathcal{V} that are distinguishable—that is, the number of subsets of \mathcal{V} for which $\rho_{\mathcal{V}}(v, v') > t$ when v and v' belong to different regions. While this construction is similar in spirit to the usual construction of packing sets in the standard reduction from testing to estimation (cf. Section 9.2.1), our bound allows us to skip the packing set construction. We can directly compute $I(V; X)$ where V takes values over the full space, as opposed to computing the mutual information $I(V'; X)$ for a random variable V' uniformly distributed over a packing set contained within \mathcal{V} . In some cases, the former calculation can be much simpler, as illustrated in examples and chapters to follow.

We now turn to providing a few consequences of Proposition 9.4.6 and Corollary 9.4.7, showing how they can be used to derive lower bounds on the minimax risk. Proposition 9.4.6 is a generalization of the classical Fano inequality (9.4.1), so it leads naturally to a generalization of the classical Fano lower bound on minimax risk, which we describe here. This reduction from estimation to testing is somewhat more general than the classical reductions, since we do not map the original estimation problem to a strict test, but rather a test that allows errors. Consider as in the standard reduction of estimation to testing in Section 9.2.1 a family of distributions $\{P_v\}_{v \in \mathcal{V}} \subset \mathcal{P}$ indexed by a finite set \mathcal{V} . This family induces an associated collection of parameters $\{\theta_v := \theta(P_v)\}_{v \in \mathcal{V}}$. Given a function $\rho_{\mathcal{V}} : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$ and a scalar t , we define the separation $\delta(t)$ of this set relative to the metric ρ on Θ via

$$\delta(t) := \sup \left\{ \delta \mid \rho(\theta_v, \theta_{v'}) \geq \delta \text{ for all } v, v' \in \mathcal{V} \text{ such that } \rho_{\mathcal{V}}(v, v') > t \right\}. \quad (9.4.11)$$

As a special case, when $t = 0$ and $\rho_{\mathcal{V}}$ is the discrete metric, this definition reduces to that of a packing set: we are guaranteed that $\rho(\theta_v, \theta_{v'}) \geq \delta(0)$ for all distinct pairs $v \neq v'$, as in the classical approach to minimax lower bounds. On the other hand, allowing for $t > 0$ lends greater flexibility to the construction, since only certain pairs θ_v and $\theta_{v'}$ are required to be well-separated.

Given a set \mathcal{V} and associated separation function (9.4.11), we assume the canonical estimation setting: nature chooses $V \in \mathcal{V}$ uniformly at random, and conditioned on this choice $V = v$, a sample X is drawn from the distribution P_v . We then have the following corollary of Proposition 9.4.6, whose argument is completely identical to that for inequality (9.2.1):

Corollary 9.4.8. *Given V uniformly distributed over \mathcal{V} with separation function $\delta(t)$, we have*

$$\mathfrak{M}_n(\theta(\mathcal{P}), \Phi \circ \rho) \geq \Phi\left(\frac{\delta(t)}{2}\right) \left[1 - \frac{I(X; V) + \log 2}{\log \frac{|\mathcal{V}|}{N_t^{\max}}} \right] \quad \text{for all } t. \quad (9.4.12)$$

Notably, using the discrete metric $\rho_{\mathcal{V}}(v, v') = \mathbf{1}\{v \neq v'\}$ and taking $t = 0$ in the lower bound (9.4.12) gives the classical Fano lower bound on the minimax risk based on constructing a packing [115, 194, 192]. We now turn to an example illustrating the use of Corollary 9.4.8 in providing a minimax lower bound on the performance of regression estimators.

Example 9.4.9 (Normal regression model): Consider the d -dimensional linear regression model $Y = X\theta + \varepsilon$, where $\varepsilon \in \mathbb{R}^n$ is i.i.d. $\mathbf{N}(0, \sigma^2)$ and $X \in \mathbb{R}^{n \times d}$ is known, but θ is not. In this case, our family of distributions is

$$\mathcal{P}_X := \left\{ Y \sim \mathbf{N}(X\theta, \sigma^2 I_{n \times n}) \mid \theta \in \mathbb{R}^d \right\} = \left\{ Y = X\theta + \varepsilon \mid \varepsilon \sim \mathbf{N}(0, \sigma^2 I_{n \times n}), \theta \in \mathbb{R}^d \right\}.$$

We then obtain the following minimax lower bound on the minimax error in squared ℓ_2 -norm: there is a universal (numerical) constant $c > 0$ such that

$$\mathfrak{M}_n(\theta(\mathcal{P}_X), \|\cdot\|_2^2) \geq c \frac{\sigma^2 d^2}{\|X\|_{\text{Fr}}^2} \geq \frac{c}{\gamma_{\max}(X/\sqrt{n})^2} \cdot \frac{\sigma^2 d}{n}, \quad (9.4.13)$$

where γ_{\max} denotes the maximum singular value. Notably, this inequality is nearly the sharpest known bound proved via Fano inequality-based methods [47], but our technique is essentially direct and straightforward.

To see inequality (9.4.13), let the set $\mathcal{V} = \{-1, 1\}^d$ be the d -dimensional hypercube, and define $\theta_v = \delta v$ for some fixed $\delta > 0$. Then letting $\rho_{\mathcal{V}}$ be the Hamming metric on \mathcal{V} and ρ be the usual ℓ_2 -norm, the associated separation function (9.4.11) satisfies $\delta(t) > \max\{\sqrt{t}, 1\}\delta$. Now, for any $t \leq \lceil d/3 \rceil$, the neighborhood size satisfies

$$N_t^{\max} = \sum_{\tau=0}^t \binom{d}{\tau} \leq 2 \binom{d}{t} \leq 2 \left(\frac{de}{t} \right)^t.$$

Consequently, for $t \leq d/6$, the ratio $|\mathcal{V}|/N_t^{\max}$ satisfies

$$\log \frac{|\mathcal{V}|}{N_t^{\max}} \geq d \log 2 - \log 2 \binom{d}{t} \geq d \log 2 - \frac{d}{6} \log(6e) - \log 2 = d \log \frac{2}{2^{1/d} \sqrt[6]{6e}} > \max \left\{ \frac{d}{6}, \log 4 \right\}$$

for $d \geq 12$. (The case $2 \leq d < 12$ can be checked directly). In particular, by taking $t = \lfloor d/6 \rfloor$ we obtain via Corollary 9.4.8 that

$$\mathfrak{M}_n(\theta(\mathcal{P}_X), \|\cdot\|_2^2) \geq \frac{\max\{\lfloor d/6 \rfloor, 2\} \delta^2}{4} \left(1 - \frac{I(Y; V) + \log 2}{\max\{d/6, 2 \log 2\}} \right).$$

But of course, for V uniform on \mathcal{V} , we have $\mathbb{E}[VV^\top] = I_{d \times d}$, and thus for V, V' independent and uniform on \mathcal{V} ,

$$\begin{aligned} I(Y; V) &\leq n \frac{1}{|\mathcal{V}|^2} \sum_{v \in \mathcal{V}} \sum_{v' \in \mathcal{V}} D_{\text{kl}}(\mathbf{N}(X\theta_v, \sigma^2 I_{n \times n}) \parallel \mathbf{N}(X\theta_{v'}, \sigma^2 I_{n \times n})) \\ &= \frac{\delta^2}{2\sigma^2} \mathbb{E} \left[\|XV - XV'\|_2^2 \right] = \frac{\delta^2}{\sigma^2} \|X\|_{\text{Fr}}^2. \end{aligned}$$

Substituting this into the preceding minimax bound, we obtain

$$\mathfrak{M}_n(\theta(\mathcal{P}_X), \|\cdot\|_2^2) \geq \frac{\max\{\lfloor d/6 \rfloor, 2\} \delta^2}{4} \left(1 - \frac{\delta^2 \|X\|_{\text{Fr}}^2 / \sigma^2 + \log 2}{\max\{d/6, 2 \log 2\}} \right).$$

Choosing $\delta^2 \asymp d\sigma^2 / \|X\|_{\text{Fr}}^2$ gives the result (9.4.13). \diamond

As a second example, we can revisit the general M-estimation results in Chapter 5.3. The construction here extends that in Example 9.4.9 to settings where there is no “true” parameter; we leave working out the details as exercises. (See Exercises 9.10 9.12, and 9.13, which applies the former two.)

Example 9.4.10 (Local minimax lower bounds for M-estimation): Let $\ell : \Theta \times \mathcal{Z} \rightarrow \mathbb{R}$ be a convex loss as in Chapter 5.3, and let it and the population loss $L(\theta) := \mathbb{E}_{P_0}[\ell(\theta, Z)]$ satisfy the conditions of Assumption A.5.1. Let $\theta_0 = \operatorname{argmin}_{\theta} L(\theta)$ have Hessian $H = \nabla^2 L(\theta_0)$ and covariance $\Sigma = \operatorname{Cov}_0(\nabla \ell(\theta_0; Z))$, so that the results in Chapter 5.3 show (roughly) that the empirical risk minimizer satisfies

$$\mathbb{E} \left[\|\hat{\theta}_n - \theta_0\|_2^2 \right] \lesssim \frac{\operatorname{tr}(H^{-1}\Sigma H^{-1})}{n} \quad \text{and} \quad \mathbb{E} \left[\|\hat{\theta}_n - \theta_0\|_{(H^{-1}\Sigma H^{-1})^{-1}}^2 \right] \lesssim \frac{d}{n},$$

where we use the notation $\|x\|_A^2 = x^\top A x$. We can show these convergence results are sharp by constructing an appropriate packing of the parameter space, induced by distributions P on \mathcal{Z} , and showing that the information the sample Z_1^n carries about the particular distribution is limited, allowing us to apply the local Fano method.

We first develop the underlying distribution family. Let $\delta \geq 0$ to be chosen. For any bounded mean-zero function $g : \mathcal{Z} \rightarrow \mathbb{R}^d$, Exercise 9.10 shows how to construct a collection of distributions $\{P_u\}_{u \in \mathbb{R}^d}$, indexed by $u \in \mathbb{R}^d$, such that

$$D_{\text{kl}}(P_u \| P_0) = \frac{1}{2} u^\top \operatorname{Cov}_0(g(X)) u + o(\|u\|^2)$$

for u near 0. Then as usual drawing a random element $V \sim \operatorname{Uniform}(\{\pm 1\}^d)$ and drawing $Z_1^n \stackrel{\text{iid}}{\sim} P_{\delta v}$ conditional on $V = v$, we have the mutual information bound

$$I(Z_1^n; V) \leq \frac{n\delta^2}{2} \operatorname{tr}(\operatorname{Cov}_0(g(Z))) + n \cdot o(\delta^2)$$

as $\delta \rightarrow 0$. For the separation of the induced parameters, let

$$\theta_v = \operatorname{argmin}_{\theta} \mathbb{E}_{P_{\delta v}}[\ell(\theta, Z)].$$

Then (see Exercise 9.13), these parameters satisfy

$$\theta_v = \theta(P_0) + \delta H^{-1} \operatorname{Cov}_0(\nabla \ell(\theta_0, Z), g(Z)) v + O(\delta^2).$$

Lastly we choose the function g . Our choice here is

$$g(Z) = (H^{-1}\Sigma H^{-1})^{-1/2} H^{-1} \nabla \ell(\theta_0, Z),$$

which is mean-zero and bounded so long as ℓ is Lipschitz near θ_0 and satisfies $\operatorname{Cov}_0(g(Z)) = \operatorname{tr}(I_d) = d$. We therefore have $I(Z_1^n; V) \lesssim n\delta^2 d$ and $\theta_v - \theta_0 = \delta(H^{-1}\Sigma H^{-1})^{-1/2} v + O(\delta^2)$. This gives a packing in the metric induced by the matrix $H\Sigma^{-1}H$, so applying Corollary 9.4.7, there exists a numerical constant $c > 0$ such that for all suitably large n we therefore obtain

$$\mathbb{P} \left(\|\hat{\theta}_n - \theta_V\|_{H\Sigma^{-1}H} \geq c\sqrt{d/n} \right) \geq \frac{1}{4}.$$

Exercises 9.12 and 9.13 work through the details. \diamond

9.5 Assouad's method

Assouad's method provides a somewhat different technique for proving lower bounds. Instead of reducing the estimation problem to a multiple hypothesis test or simpler estimation problem, as with Le Cam's method and Fano's method from the preceding lectures, here we transform the original estimation problem into multiple binary hypothesis testing problems, using the structure of the problem in an essential way. Assouad's method applies only problems where the loss we care about is naturally related to identification of individual points on a hypercube. In simple or standard problems, Assouad's method rarely provides stronger lower bounds than the local Fano methods; its true power lies in its applications to *adaptive* problems, where one may choose points at which to query a statistical model. We develop this idea in the exercises (see Exercise 9.6), and also leverage it in Chapter 18, Section 18.5.2 to prove lower bounds for bandit problems, where one must balance adaptive exploration of a function with performance.

9.5.1 Well-separated problems

To describe the method, we begin by encoding a notion of separation and loss, similar to what we did in the classical reduction of estimation to testing. For some $d \in \mathbb{N}$, let $\mathcal{V} = \{-1, 1\}^d$, and let us consider a family $\{P_v\}_{v \in \mathcal{V}} \subset \mathcal{P}$ indexed by the hypercube. We say that the family P_v induces a 2δ -Hamming separation for the loss $\Phi \circ \rho$ if there exists a function $\hat{v} : \theta(\mathcal{P}) \rightarrow \{-1, 1\}^d$ satisfying

$$\Phi(\rho(\theta, \theta(P_v))) \geq 2\delta \sum_{j=1}^d \mathbf{1}\{\hat{v}(\theta)_j \neq v_j\}. \quad (9.5.1)$$

That is, we can take the parameter θ and test the individual indices via \hat{v} .

Example 9.5.1 (Estimation in ℓ_1 -error): Suppose we have a family of multivariate Laplace distributions on \mathbb{R}^d —distributions with density proportional to $p(x) \propto \exp(-\|x - \mu\|_1)$ —and we wish to estimate the mean in ℓ_1 -distance. For $v \in \{-1, 1\}^d$ and some fixed $\delta > 0$ let p_v be the density

$$p_v(x) = \frac{1}{2} \exp(-\|x - \delta v\|_1),$$

which has mean $\theta(P_v) = \delta v$. Under the ℓ_1 -loss, we have for any $\theta \in \mathbb{R}^d$ that

$$\|\theta - \theta(P_v)\|_1 = \sum_{j=1}^d |\theta_j - \delta v_j| \geq \delta \sum_{j=1}^d \mathbf{1}\{\text{sign}(\theta_j) \neq v_j\},$$

so that this family induces a δ -Hamming separation for the ℓ_1 -loss. \diamond

9.5.2 From estimation to multiple binary tests

As in the standard reduction from estimation to testing, we consider the following random process: nature chooses a vector $V \in \{-1, 1\}^d$ uniformly at random, after which the sample X is drawn from the distribution P_v conditional on $V = v$. Then, if we let $\mathbb{P}_{\pm j}$ denote the joint distribution over the random index V and X conditional on the j th coordinate $V_j = \pm 1$, we obtain the following sharper version of Assouad's lemma [10] (see also the paper [8]); we provide a proof in Section 9.6.3 to follow.

Lemma 9.5.2. *Under the conditions of the previous paragraph, we have*

$$\mathfrak{M}(\theta(\mathcal{P}), \Phi \circ \rho) \geq \delta \sum_{j=1}^d \inf_{\Psi} [\mathbb{P}_{+j}(\Psi(X) \neq +1) + \mathbb{P}_{-j}(\Psi(X) \neq -1)].$$

While Lemma 9.5.2 requires conditions on the loss Φ and metric ρ for the separation condition (9.5.1) to hold, it is sometimes easier to apply than Fano's method. Moreover, while we will not address this in class, several researchers [8, 74] have noted that it appears to allow easier application in so-called “interactive” settings—those for which the sampling of the X_i may not be precisely i.i.d. It is closely related to Le Cam's method, discussed previously, as we see that if we define $P_{+j} = 2^{1-d} \sum_{v: v_j=1} P_v$ (and similarly for $-j$), Lemma 9.5.2 is equivalent to

$$\mathfrak{M}(\theta(\mathcal{P}), \Phi \circ \rho) \geq \delta \sum_{j=1}^d [1 - \|P_{+j} - P_{-j}\|_{\text{TV}}]. \quad (9.5.2)$$

There are standard weakenings of the lower bound (9.5.2) (and Lemma 9.5.2). We give one such weakening. First, we note that the total variation is convex, so that if we define $P_{v,+j}$ to be the distribution P_v where coordinate j takes the value $v_j = 1$ (and similarly for $P_{v,-j}$), we have

$$P_{+j} = \frac{1}{2^d} \sum_{v \in \{-1,1\}^d} P_{v,+j} \quad \text{and} \quad P_{-j} = \frac{1}{2^d} \sum_{v \in \{-1,1\}^d} P_{v,-j}.$$

Thus, by the triangle inequality, we have

$$\begin{aligned} \|P_{+j} - P_{-j}\|_{\text{TV}} &= \left\| \frac{1}{2^d} \sum_{v \in \{-1,1\}^d} P_{v,+j} - \frac{1}{2^d} \sum_{v \in \{-1,1\}^d} P_{v,-j} \right\|_{\text{TV}} \\ &\leq \frac{1}{2^d} \sum_{v \in \{-1,1\}^d} \|P_{v,+j} - P_{v,-j}\|_{\text{TV}} \leq \max_{v,j} \|P_{v,+j} - P_{v,-j}\|_{\text{TV}}. \end{aligned}$$

Then as long as the loss satisfies the per-coordinate separation (9.5.1), we obtain the following:

$$\mathfrak{M}(\theta(\mathcal{P}), \Phi \circ \rho) \geq d\delta \left(1 - \max_{v,j} \|P_{v,+j} - P_{v,-j}\|_{\text{TV}} \right). \quad (9.5.3)$$

This most common version of Assouad's lemma sometimes too brutally controls $\|P_{+j} - P_{-j}\|_{\text{TV}}$.

We also note that by the Cauchy-Schwarz inequality and convexity of the variation-distance, we have

$$\sum_{j=1}^d \|P_{+j} - P_{-j}\|_{\text{TV}} \leq \sqrt{d} \left(\sum_{j=1}^d \|P_{+j} - P_{-j}\|_{\text{TV}}^2 \right)^{1/2} \leq \sqrt{d} \left(\sum_{j=1}^d \frac{1}{2^d} \sum_v \|P_{v,+j} - P_{v,-j}\|_{\text{TV}}^2 \right)^{1/2},$$

and consequently we have a not quite so terribly weak version of inequality (9.5.2):

$$\mathfrak{M}(\theta(\mathcal{P}), \Phi \circ \rho) \geq \delta d \left[1 - \left(\frac{1}{2^d d} \sum_{j=1}^d \sum_{v \in \{-1,1\}^d} \|P_{v,+j} - P_{v,-j}\|_{\text{TV}}^2 \right)^{1/2} \right]. \quad (9.5.4)$$

Regardless of whether we use the sharper version (9.5.2) or weakened versions (9.5.3) or (9.5.4), the technique is essentially the same. We seek a setting of the distributions P_v so that the probability of making a mistake in the hypothesis test of Lemma 9.5.2 is high enough—say $1/2$ —or the variation distance is small enough—such as $\|P_{+j} - P_{-j}\|_{\text{TV}} \leq 1/2$ for all j . Once this is satisfied, we obtain a minimax lower bound of the form

$$\mathfrak{M}(\theta(\mathcal{P}), \Phi \circ \rho) \geq \delta \sum_{j=1}^d \left[1 - \frac{1}{2}\right] = \frac{d\delta}{2}.$$

9.5.3 Example applications of Assouad's method

We now provide two example applications of Assouad's method. The first is a standard finite-dimensional lower bound, where we provide a lower bound in a normal mean estimation problem. For the second, we consider estimation in a logistic regression problem, showing a similar lower bound. In Section 10.1 to follow, we show how to use Assouad's method to prove strong lower bounds in a standard nonparametric problem.

Example 9.5.3 (Normal mean estimation): For some $\sigma^2 > 0$ and $d \in \mathbb{N}$, we consider estimation of mean parameter for the normal location family

$$\mathcal{N} := \left\{ \mathbf{N}(\theta, \sigma^2 I_{d \times d}) : \theta \in \mathbb{R}^d \right\}$$

in squared Euclidean distance. We now show how for this family, the sharp Assouad's method implies the lower bound

$$\mathfrak{M}_n(\theta(\mathcal{N}), \|\cdot\|_2^2) \geq \frac{d\sigma^2}{8n}. \quad (9.5.5)$$

Up to constant factors, this bound is sharp; the sample mean has mean squared error $d\sigma^2/n$.

We proceed in (essentially) the usual way we have set up. Fix some $\delta > 0$ and define $\theta_v = \delta v$, taking $P_v = \mathbf{N}(\theta_v, \sigma^2 I_{d \times d})$ to be the normal distribution with mean θ_v . In this case, we see that the hypercube structure is natural, as our loss function decomposes on coordinates: we have $\|\theta - \theta_v\|_2^2 \geq \delta^2 \sum_{j=1}^d \mathbf{1}\{\text{sign}(\theta_j) \neq v_j\}$. The family P_v thus induces a δ^2 -Hamming separation for the loss $\|\cdot\|_2^2$, and by Assouad's method (9.5.2), we have

$$\mathfrak{M}_n(\theta(\mathcal{N}), \|\cdot\|_2^2) \geq \frac{\delta^2}{2} \sum_{j=1}^d \left[1 - \|P_{+j}^n - P_{-j}^n\|_{\text{TV}}\right],$$

where $P_{\pm j}^n = 2^{1-d} \sum_{v: v_j = \pm 1} P_v^n$. It remains to provide upper bounds on $\|P_{+j}^n - P_{-j}^n\|_{\text{TV}}$. By the convexity of $\|\cdot\|_{\text{TV}}^2$ and Pinsker's inequality, we have

$$\|P_{+j}^n - P_{-j}^n\|_{\text{TV}}^2 \leq \max_{d_{\text{ham}}(v, v') \leq 1} \|P_v^n - P_{v'}^n\|_{\text{TV}}^2 \leq \frac{1}{2} \max_{d_{\text{ham}}(v, v') \leq 1} D_{\text{kl}}(P_v^n \| P_{v'}^n).$$

But of course, for any v and v' differing in only 1 coordinate,

$$D_{\text{kl}}(P_v^n \| P_{v'}^n) = \frac{n}{2\sigma^2} \|\theta_v - \theta_{v'}\|_2^2 = \frac{2n}{\sigma^2} \delta^2,$$

giving the minimax lower bound

$$\mathfrak{M}_n(\theta(\mathcal{N}), \|\cdot\|_2^2) \geq 2\delta^2 \sum_{j=1}^d \left[1 - \sqrt{2n\delta^2/\sigma^2}\right].$$

Choosing $\delta^2 = \sigma^2/8n$ gives the claimed lower bound (9.5.5). \diamond

Example 9.5.4 (Logistic regression): In this example, consider the logistic regression model, where we have known (fixed) regressors $X_i \in \mathbb{R}^d$ and an unknown parameter $\theta \in \mathbb{R}^d$; the goal is to estimate θ after observing a sequence of $Y_i \in \{-1, 1\}$, where for $y \in \{-1, 1\}$ we have

$$P(Y_i = y \mid X_i, \theta) = \frac{1}{1 + \exp(-yX_i^\top \theta)}.$$

Denote this family by \mathcal{P}_{\log} , and for $P \in \mathcal{P}_{\log}$, let $\theta(P)$ be the predictor vector θ . We would like to estimate the vector θ in squared ℓ_2 error. As in Example 9.5.3, if we choose some $\delta > 0$ and for each $v \in \{-1, 1\}^d$, we set $\theta_v = \delta v$, then we have the δ^2 -separation in Hamming metric $\|\theta - \theta_v\|_2^2 \geq \delta^2 \sum_{j=1}^d \mathbf{1}\{\text{sign}(\theta_j) \neq v_j\}$. Let P_v^n denote the distribution of the n independent observations Y_i when $\theta = \theta_v$. Then we have by Assouad's lemma (and the weakening (9.5.4)) that

$$\begin{aligned} \mathfrak{M}_n(\theta(\mathcal{P}_{\log}), \|\cdot\|_2^2) &\geq \frac{\delta^2}{2} \sum_{j=1}^d \left[1 - \|P_{+j}^n - P_{-j}^n\|_{\text{TV}}\right] \\ &\geq \frac{d\delta^2}{2} \left[1 - \left(\frac{1}{d} \sum_{j=1}^d \frac{1}{2^d} \sum_{v \in \{-1, 1\}^d} \|P_{v,+j}^n - P_{v,-j}^n\|_{\text{TV}}^2\right)^{\frac{1}{2}}\right]. \end{aligned} \quad (9.5.6)$$

It remains to bound $\|P_{v,+j}^n - P_{v,-j}^n\|_{\text{TV}}^2$ to find our desired lower bound. To that end, use the shorthands $p_v(x) = 1/(1 + \exp(\delta x^\top v))$ and let $D_{\text{kl}}(p\|q)$ be the binary KL-divergence between Bernoulli(p) and Bernoulli(q) distributions. Then Pinsker's inequality (recall Proposition 2.2.8) implies that for any v, v' ,

$$\begin{aligned} \|P_v^n - P_{v'}^n\|_{\text{TV}} &\leq \frac{1}{4} [D_{\text{kl}}(P_v^n \| P_{v'}^n) + D_{\text{kl}}(P_{v'}^n \| P_v^n)] \\ &= \frac{1}{4} \sum_{i=1}^n [D_{\text{kl}}(p_v(X_i) \| p_{v'}(X_i)) + D_{\text{kl}}(p_{v'}(X_i) \| p_v(X_i))]. \end{aligned}$$

Let us upper bound the final KL-divergence. Let $p_a = 1/(1 + e^a)$ and $p_b = 1/(1 + e^b)$. We claim that

$$D_{\text{kl}}(p_a \| p_b) + D_{\text{kl}}(p_b \| p_a) \leq (a - b)^2. \quad (9.5.7)$$

Deferring the proof of claim (9.5.7), we immediately see that

$$\|P_v^n - P_{v'}^n\|_{\text{TV}} \leq \frac{\delta^2}{4} \sum_{i=1}^n \left(X_i^\top (v - v')\right)^2.$$

Now we recall inequality (9.5.6) for motivation, and we see that the preceding display implies

$$\frac{1}{2^d d} \sum_{j=1}^d \sum_{v \in \{-1, 1\}^d} \|P_{v,+j}^n - P_{v,-j}^n\|_{\text{TV}}^2 \leq \frac{\delta^2}{4d} \frac{1}{2^d} \sum_{v \in \{-1, 1\}^d} \sum_{j=1}^d \sum_{i=1}^n (2X_{ij})^2 = \frac{\delta^2}{d} \sum_{i=1}^n \sum_{j=1}^d X_{ij}^2.$$

Replacing the final double sum with $\|X\|_{\text{Fr}}^2$, where X is the matrix of the X_i , we have

$$\mathfrak{M}_n(\theta(\mathcal{P}_{\log}), \|\cdot\|_2^2) \geq \frac{d\delta^2}{2} \left[1 - \left(\frac{\delta^2}{d} \|X\|_{\text{Fr}}^2 \right)^{\frac{1}{2}} \right].$$

Setting $\delta^2 = d/4 \|X\|_{\text{Fr}}^2$, we obtain

$$\mathfrak{M}_n(\theta(\mathcal{P}_{\log}), \|\cdot\|_2^2) \geq \frac{d\delta^2}{4} = \frac{d^2}{16 \|X\|_{\text{Fr}}^2} = \frac{d}{n} \cdot \frac{1}{16 \frac{1}{dn} \sum_{i=1}^n \|X_i\|_2^2}.$$

That is, we have a minimax lower bound scaling roughly as d/n for logistic regression, where “large” X_i (in ℓ_2 -norm) suggest that we may obtain better performance in estimation. This is intuitive, as a larger X_i gives a better signal to noise ratio.

We return to prove the claim (9.5.7). Indeed, by a straightforward expansion, we have

$$\begin{aligned} D_{\text{kl}}(p_a \| p_b) + D_{\text{kl}}(p_b \| p_a) &= p_a \log \frac{p_a}{p_b} + (1 - p_a) \log \frac{1 - p_a}{1 - p_b} + p_b \log \frac{p_b}{p_a} + (1 - p_b) \log \frac{1 - p_b}{1 - p_a} \\ &= (p_a - p_b) \log \frac{p_a}{p_b} + (p_b - p_a) \log \frac{1 - p_a}{1 - p_b} = (p_a - p_b) \log \left(\frac{p_a}{1 - p_a} \frac{1 - p_b}{p_b} \right). \end{aligned}$$

Now note that $p_a/(1 - p_a) = e^{-a}$ and $(1 - p_b)/p_b = e^b$. Thus we obtain

$$D_{\text{kl}}(p_a \| p_b) + D_{\text{kl}}(p_b \| p_a) = \left(\frac{1}{1 + e^a} - \frac{1}{1 + e^b} \right) \log(e^{b-a}) = (b - a) \left(\frac{1}{1 + e^a} - \frac{1}{1 + e^b} \right)$$

Assume without loss of generality that $b \geq a$. Noting that $e^x \geq 1 + x$ by convexity, we have

$$\frac{1}{1 + e^a} - \frac{1}{1 + e^b} = \frac{e^b - e^a}{(1 + e^a)(1 + e^b)} \leq \frac{e^b - e^a}{e^b} = 1 - e^{a-b} \leq 1 - (1 + (a - b)) = b - a,$$

yielding claim (9.5.7). \diamond

9.6 Deferred proofs

9.6.1 Proof of Proposition 9.4.6

Our argument for proving the proposition parallels that of the classical Fano inequality by Cover and Thomas [57]. Letting E be a $\{0, 1\}$ -valued indicator variable for the event $\rho(\hat{V}, V) \leq t$, we compute the entropy $H(E, V \mid \hat{V})$ in two different ways. On one hand, by the chain rule for entropy, we have

$$H(E, V \mid \hat{V}) = H(V \mid \hat{V}) + \underbrace{H(E \mid V, \hat{V})}_{=0}, \quad (9.6.1)$$

where the final term vanishes since E is (V, \hat{V}) -measurable. On the other hand, we also have

$$H(E, V \mid \hat{V}) = H(E \mid \hat{V}) + H(V \mid E, \hat{V}) \leq H(E) + H(V \mid E, \hat{V}),$$

using the fact that conditioning reduces entropy. Applying the definition of conditional entropy yields

$$H(V \mid E, \hat{V}) = \mathbb{P}(E = 0)H(V \mid E = 0, \hat{V}) + \mathbb{P}(E = 1)H(V \mid E = 1, \hat{V}),$$

and we upper bound each of these terms separately. For the first term, we have

$$H(V \mid E = 0, \hat{V}) \leq \log(|\mathcal{V}| - N_t^{\min}),$$

since conditioned on the event $E = 0$, the random variable V may take values in a set of size at most $|\mathcal{V}| - N_t^{\min}$. For the second, we have

$$H(V \mid E = 1, \hat{V}) \leq \log N_t^{\max},$$

since conditioned on $E = 1$, or equivalently on the event that $\rho(\hat{V}, V) \leq t$, we are guaranteed that V belongs to a set of cardinality at most N_t^{\max} .

Combining the pieces and noting $\mathbb{P}(E = 0) = P_t$, we have proved that

$$H(E, V \mid \hat{V}) \leq H(E) + P_t \log(|\mathcal{V}| - N_t^{\min}) + (1 - P_t) \log N_t^{\max}.$$

Combining this inequality with our earlier equality (9.6.1), we see that

$$H(V \mid \hat{V}) \leq H(E) + P_t \log(|\mathcal{V}| - N_t^{\min}) + (1 - P_t) \log N_t^{\max}.$$

Since $H(E) = h_2(P_t)$, the claim (9.4.9) follows.

9.6.2 Proof of Corollary 9.4.7

First, by the information-processing inequality [e.g. 57, Chapter 2], we have $I(V; \hat{V}) \leq I(V; X)$, and hence $H(V \mid X) \leq H(V \mid \hat{V})$. Since $h_2(P_t) \leq \log 2$, inequality (9.4.9) implies that

$$H(V \mid X) - \log N_t^{\max} \leq H(V \mid \hat{V}) - \log N_t^{\max} \leq \mathbb{P}(\rho(\hat{V}, V) > t) \log \frac{|\mathcal{V}| - N_t^{\min}}{N_t^{\max}} + \log 2.$$

Rearranging the preceding equations yields

$$\mathbb{P}(\rho(\hat{V}, V) > t) \geq \frac{H(V \mid X) - \log N_t^{\max} - \log 2}{\log \frac{|\mathcal{V}| - N_t^{\min}}{N_t^{\max}}}. \quad (9.6.2)$$

Note that this bound holds without any assumptions on the distribution of V .

By definition, we have $I(V; X) = H(V) - H(V \mid X)$. When V is uniform on \mathcal{V} , we have $H(V) = \log |\mathcal{V}|$, and hence $H(V \mid X) = \log |\mathcal{V}| - I(V; X)$. Substituting this relation into the bound (9.6.2) yields the inequality

$$\mathbb{P}(\rho(\hat{V}, V) > t) \geq \frac{\log \frac{|\mathcal{V}|}{N_t^{\max}}}{\log \frac{|\mathcal{V}| - N_t^{\min}}{N_t^{\max}}} - \frac{I(V; X) + \log 2}{\log \frac{|\mathcal{V}| - N_t^{\min}}{N_t^{\max}}} \geq 1 - \frac{I(V; X) + \log 2}{\log \frac{|\mathcal{V}|}{N_t^{\max}}}.$$

9.6.3 Proof of Lemma 9.5.2

Fix an (arbitrary) estimator $\hat{\theta}$. By assumption (9.5.1), we have

$$\Phi(\rho(\theta, \theta(P_v))) \geq 2\delta \sum_{j=1}^d \mathbf{1}\{[\hat{v}(\theta)]_j \neq v_j\}.$$

Taking expectations, we see that

$$\begin{aligned} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\Phi(\rho(\hat{\theta}(X), \theta(P))) \right] &\geq \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \mathbb{E}_{P_v} \left[\Phi(\rho(\hat{\theta}(X), \theta_v)) \right] \\ &\geq \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} 2\delta \sum_{j=1}^d \mathbb{E}_{P_v} \left[\mathbf{1} \left\{ [\psi(\hat{\theta})]_j \neq v_j \right\} \right] \end{aligned}$$

as the average is smaller than the maximum of a set and using the separation assumption (9.5.1). Recalling the definition of the mixtures $\mathbb{P}_{\pm j}$ as the joint distribution of V and X conditional on $V_j = \pm 1$, we swap the summation orders to see that

$$\begin{aligned} \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P_v \left([\hat{v}(\hat{\theta})]_j \neq v_j \right) &= \frac{1}{|\mathcal{V}|} \sum_{v: v_j=1} P_v \left([\hat{v}(\hat{\theta})]_j \neq v_j \right) + \frac{1}{|\mathcal{V}|} \sum_{v: v_j=-1} P_v \left([\hat{v}(\hat{\theta})]_j \neq v_j \right) \\ &= \frac{1}{2} \mathbb{P}_{+j} \left([\hat{v}(\hat{\theta})]_j \neq v_j \right) + \frac{1}{2} \mathbb{P}_{-j} \left([\hat{v}(\hat{\theta})]_j \neq v_j \right). \end{aligned}$$

This gives the statement claimed in the lemma, while taking an infimum over all testing procedures $\Psi : \mathcal{X} \rightarrow \{-1, +1\}$ gives the claim (9.5.2).

9.7 Bibliography

For a fuller technical introduction into nonparametric estimation, see the book by Tsybakov [182]. Has'minskii [109].

The material in Section 10.2 is based on a paper of Yang and Barron [192].

9.8 Exercises

Exercise 9.1 (A generalized version of Fano's inequality; cf. Proposition 9.4.6): Let \mathcal{V} and $\hat{\mathcal{V}}$ be arbitrary sets, and suppose that π is a (prior) probability measure on \mathcal{V} , where V is distributed according to π . Let $V \rightarrow X \rightarrow \hat{V}$ be Markov chain, where V takes values in \mathcal{V} and \hat{V} takes values in $\hat{\mathcal{V}}$. Let $\mathcal{N} \subset \mathcal{V} \times \hat{\mathcal{V}}$ denote a measurable subset of $\mathcal{V} \times \hat{\mathcal{V}}$ (a collection of neighborhoods), and for any $\hat{v} \in \hat{\mathcal{V}}$, denote the slice

$$\mathcal{N}_{\hat{v}} := \{v \in \mathcal{V} : (v, \hat{v}) \in \mathcal{N}\}. \quad (9.8.1)$$

That is, \mathcal{N} denotes the neighborhoods of points v for which we do not consider a prediction \hat{v} for v to be an error, and the slices (9.8.1) index the neighborhoods. Define the “volume” constants

$$p^{\max} := \sup_{\hat{v}} \pi(V \in \mathcal{N}_{\hat{v}}) \quad \text{and} \quad p^{\min} := \inf_{\hat{v}} \pi(V \in \mathcal{N}_{\hat{v}}).$$

Define the error probability $P_{\text{error}} = \mathbb{P}[(V, \hat{V}) \notin \mathcal{N}]$ and entropy $h_2(p) = -p \log p - (1-p) \log(1-p)$.

(a) Prove that for any Markov chain $V \rightarrow X \rightarrow \hat{V}$, we have

$$h_2(P_{\text{error}}) + P_{\text{error}} \log \frac{1 - p^{\min}}{p^{\max}} \geq \log \frac{1}{p^{\max}} - I(V; \hat{V}). \quad (9.8.2)$$

(b) Conclude from inequality (9.8.2) that

$$\mathbb{P}[(V, \hat{V}) \notin \mathcal{N}] \geq 1 - \frac{I(V; X) + \log 2}{\inf_{\hat{v}} \log \frac{1}{\pi(\mathcal{N}_{\hat{v}})}}.$$

(c) Now we give a version explicitly using distances. Let $\mathcal{V} \subset \mathbb{R}^d$ and define $\mathcal{N} = \{(v, v') : \|v - v'\| \leq \delta\}$ to be the points within δ of one another. Let \mathbb{B}_v denote the $\|\cdot\|$ -ball of radius 1 centered at v . Conclude that for any prior π on \mathbb{R}^d that

$$\mathbb{P}(\|V - \hat{V}\|_2 \geq \delta) \geq 1 - \frac{I(V; X) + \log 2}{\log \frac{1}{\sup_v \pi(\delta \mathbb{B}_v)}}.$$

Exercise 9.2: In this question, we will show that the minimax rate of estimation for the parameter of a uniform distribution (in squared error) scales as $1/n^2$. In particular, assume that $X_i \stackrel{\text{iid}}{\sim} \text{Uniform}(\theta, \theta + 1)$, meaning that X_i have densities $p(x) = \mathbf{1}\{x \in [\theta, \theta + 1]\}$. Let $X_{(1)} = \min_i \{X_i\}$ denote the first order statistic.

(a) Prove that

$$\mathbb{E}[(X_{(1)} - \theta)^2] = \frac{2}{(n+1)(n+2)}.$$

(Hint: the fact that $\mathbb{E}[Z] = \int_0^\infty \mathbb{P}(Z \geq t) dt$ for any positive Z may be useful.)

(b) Using Le Cam's two-point method, show that the minimax rate for estimation of $\theta \in \mathbb{R}$ for the uniform family $\mathcal{U} = \{\text{Uniform}(\theta, \theta + 1) : \theta \in \mathbb{R}\}$ in squared error has lower bound c/n^2 , where c is a numerical constant.

Exercise 9.3 (Sign identification in sparse linear regression): In sparse linear regression, we have n observations $Y_i = \langle X_i, \theta^* \rangle + \varepsilon_i$, where $X_i \in \mathbb{R}^d$ are known (fixed) matrices and the vector θ^* has a small number $k \ll d$ of non-zero indices, and $\varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$. In this problem, we investigate the problem of *sign recovery*, that is, identifying the vector of signs $\text{sign}(\theta_j^*)$ for $j = 1, \dots, d$, where $\text{sign}(0) = 0$.

Assume we have the following process: fix a signal threshold $\theta_{\min} > 0$. First, a vector $S \in \{-1, 0, 1\}^d$ is chosen uniformly at random from the set of vectors $\mathcal{S}_k := \{s \in \{-1, 0, 1\}^d : \|s\|_1 = k\}$. Then we define vectors θ^s so that $\theta_j^s = \theta_{\min} s_j$, and conditional on $S = s$, we observe

$$Y = X\theta^s + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_{n \times n}).$$

(Here $X \in \mathbb{R}^{n \times d}$ is a known fixed matrix.)

(a) Use Fano's inequality to show that for any estimator \hat{S} of S , we have

$$\mathbb{P}(\hat{S} \neq S) \geq \frac{1}{2} \quad \text{unless} \quad n \geq c \frac{\frac{d}{k} \log \binom{d}{k}}{\|n^{-1/2} X\|_{\text{Fr}}^2} \frac{\sigma^2}{\theta_{\min}^2},$$

where c is a numerical constant. You may assume that $k \geq 4$ or $\log \binom{d}{k} \geq 4 \log 2$.

(b) Assume that $X \in \{-1, 1\}^{n \times d}$. Give a lower bound on how large n must be for sign recovery. Give a one sentence interpretation of σ^2/θ_{\min}^2 .

Exercise 9.4 (Multiple hypothesis testing and recovery): A p -value for a null hypothesis H is any random variable $Y \in [0, 1]$ such that $\mathbb{P}(Y \leq u) \leq u$ for all $u \in [0, 1]$ whenever H holds, so that it is no more likely to be small than a $\text{Uniform}[0, 1]$ random variable (as Y being quite small is evidence against the hypothesis). In the multiple hypothesis testing problem, we consider n null hypotheses H_1, \dots, H_n , and an associated p -value $Y_i \in [0, 1]$ for each. The goal in multiple hypothesis testing is to reject as many hypotheses as possible without making false discoveries, that is, rejecting an index i for which the null H_i holds.

We can model this via an estimation problem of estimating a binary vector $v \in \{0, 1\}^n$ indicating which hypotheses are null ($v_i = 0$) and which are non-null. To mathematize this, let $V \in \{0, 1\}^n$ have i.i.d. entries $V_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\epsilon)$, and conditional on $V_i = 0$, draw $Y_i \sim P_0 := \text{Uniform}([0, 1])$, and conditional on $V_i = 1$, draw $Y_i \stackrel{\text{iid}}{\sim} P_1$, where P_1 has support $[0, 1]$; we consider the consequences of making various choices for P_1 .

(a) Let $P_1 = \text{Uniform}[0, \tau]$ for some $\tau < 1$. Show that $\inf_{\hat{v}} \mathbb{P}(V_i \neq \hat{v}(Y_1^n)) = \min\{\tau(1 - \epsilon), \epsilon\}$.

(b) For a fixed value $0 < \tau < 1$, let

$$P_1 = (1 - \tau)\text{Uniform}[0, \tau] + \tau\text{Uniform}[\tau, 1],$$

so that $Y \sim P_1$ has density $p_1(y) = \frac{1-\tau}{\tau} \mathbf{1}\{0 \leq y \leq \tau\} + \frac{\tau}{1-\tau} \mathbf{1}\{\tau < y \leq 1\}$. Show that

$$\inf_{\hat{v}} \mathbb{P}(V_i \neq \hat{v}(Y_1^n)) = \min\{\tau(1 - \epsilon), \epsilon(1 - \tau)\} + \min\{(1 - \tau)(1 - \epsilon), \epsilon\tau\}.$$

(c) Let the setting of part (b) hold. $\mathcal{I}_0 := \{i \in [n] \mid V_i = 0\}$ and $\mathcal{I}_1 := \{i \in [n] \mid V_i = 1\}$ be the collections of null and non-null indices, respectively. Assume the *sparse testing* regime, where we take

$$\epsilon = n^{-\beta} \quad \text{and} \quad \tau = n^{-r},$$

where $\beta \in (\frac{1}{2}, 1)$ and $r \in (0, \infty)$. For a test $\hat{v} : [0, 1]^n \rightarrow \{0, 1\}^n$, let $R := \{i \mid \hat{v}_i = 1\}$ be the set of rejected null hypotheses. Show that if $r < \beta$,

$$\mathbb{E} \left[\frac{\text{card}(R \cap \mathcal{I}_0) + \text{card}(R^c \cap \mathcal{I}_1)}{\text{card}(\mathcal{I}_1) + 1} \right] \geq 1 - o(1).$$

That is, the number of mistakes scales at least as the number of true non-nulls.

Exercise 9.5 (Benjamini-Hochberg procedures and optimal recovery): Consider a multiple hypothesis testing problem in which a statistician receives a set of p -values, which we denote $\{Y_1, \dots, Y_n\}$, and seeks to test each of n null hypotheses H_1, \dots, H_n ; under the null H_i we have $\mathbb{P}_{H_i}(Y_i \leq u) \leq u$ for all $u \in [0, 1]$. Fix a desired level $\alpha \in (0, 1)$. The *Bonferroni correction* (that is, the union bound) rejects those hypothesis satisfying

$$Y_i \leq \frac{\alpha}{n}.$$

Immediately, we obtain $\mathbb{P}(\text{false rejection}) \leq \alpha$. In contrast, the Benjamini-Hochberg (BH) procedure [23] sorts the values $Y_{(1)} \leq Y_{(2)} \leq \dots \leq Y_{(n)}$, finds the largest value

$$\hat{k} := \max \left\{ k \in \mathbb{N} \mid Y_{(k)} \leq \frac{k\alpha}{n} \right\}, \quad (9.8.3)$$

then rejects $H_{(1)}, \dots, H_{(\widehat{k})}$, the associated nulls (where $\widehat{k} = 0$ if $Y_{(k)} > \frac{k\alpha}{n}$ for each k). [Benjamini and Hochberg](#) prove that if F is the number of *false discoveries*, meaning hypotheses rejected that were null, then so long as the nulls are independent

$$\mathbb{E} \left[F / \max\{\widehat{k}, 1\} \right] \leq \alpha.$$

We compare these two procedures under the distributional setting of Exercise 9.4 part (b), where $\tau = n^{-r}$ and $\epsilon = n^{-\beta}$ for $r > \beta$ and $\beta \in (\frac{1}{2}, 1)$. Let $\mathcal{I}_0 = \{i \in [n] \mid V_i = 0\}$ and $\mathcal{I}_1 = \{i \in [n] \mid V_i = 1\}$ be the null and non-null hypotheses, respectively.

- (a) Assume $r < 1$ and let $R_{\text{BC}} = \{i \in [n] \mid Y_i \leq \frac{\alpha}{n}\}$ be the hypotheses the Bonferroni correction rejects. Show that

$$\mathbb{E}[\text{card}(R_{\text{BC}}^c \cap \mathcal{I}_1)] \geq \epsilon n(1 - o(1)) = n^{1-\beta}(1 - o(1)) \quad \text{and} \quad \mathbb{E} \left[\frac{\text{card}(R_{\text{BC}}^c \cap \mathcal{I}_1)}{\text{card}(\mathcal{I}_1) + 1} \right] \geq 1 - o(1).$$

For the remainder of the question, let $R \subset [n]$ be the indices the BH procedure (9.8.3) rejects.

- (b) Show that if $r > \beta$, then for each $\delta > 0$,

$$\mathbb{P}(\text{card}(R \cap \mathcal{I}_1) \leq (1 - \delta) \text{card}(\mathcal{I}_1)) \rightarrow 0.$$

- (c) Show that if $r > \beta$, then for any sequence $\alpha_n \rightarrow 0$ slowly enough in the BH procedure (9.8.3),

$$\mathbb{E} \left[\frac{\text{card}(R \cap \mathcal{I}_0) + \text{card}(R^c \cap \mathcal{I}_1)}{\text{card}(\mathcal{I}_1) + 1} \right] \rightarrow 0.$$

- (d) Comparing with exercise 9.4, what can we say about the relative merits of the Bonferroni correction and the BH procedure (9.8.3)?

See also the paper [1] for more on optimality of procedures of the form (9.8.3).

Exercise 9.6: In this question, we study the question of whether adaptivity can give better estimation performance for linear regression problems. That is, for $i = 1, \dots, n$, assume that we observe variables Y_i in the usual linear regression setup,

$$Y_i = \langle X_i, \theta \rangle + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathbf{N}(0, \sigma^2), \quad (9.8.4)$$

where $\theta \in \mathbb{R}^d$ is unknown. But now, based on observing $Y_1^{i-1} = \{Y_1, \dots, Y_{i-1}\}$, we allow an adaptive choice of the next predictor variables $X_i \in \mathbb{R}^d$. Let $\mathcal{L}_{\text{ada}}^n(\mathbf{F}^2)$ denote the family of linear regression problems under this adaptive setting (with n observations) where we constrain the Frobenius norm of the data matrix $X^\top = [X_1 \cdots X_n]$, $X \in \mathbb{R}^{n \times d}$, to have bound $\|X\|_{\text{Fr}}^2 = \sum_{i=1}^n \|X_i\|_2^2 \leq \mathbf{F}^2$. We use Assouad's method to show that the minimax mean-squared error satisfies the following bound:

$$\mathfrak{M}(\mathcal{L}_{\text{ada}}^n(\mathbf{F}^2), \|\cdot\|_2^2) := \inf_{\widehat{\theta}} \sup_{\theta \in \mathbb{R}^d} \mathbb{E}[\|\widehat{\theta} - \theta\|_2^2] \geq \frac{d\sigma^2}{n} \cdot \frac{1}{16 \frac{1}{dn} \mathbf{F}^2}. \quad (9.8.5)$$

Here the infimum is taken over all adaptive procedures satisfying $\|X\|_{\text{Fr}}^2 \leq \mathbf{F}^2$.

In general, when we choose X_i based on the observations Y_1^{i-1} , we are taking $X_i = F_i(Y_1^{i-1}, U_1^i)$, where U_i is a random variable independent of ε_i and Y_1^{i-1} and F_i is some function. Justify the following steps in the proof of inequality (9.8.5):

- (i) Assume that nature chooses $v \in \mathcal{V} = \{-1, 1\}^d$ uniformly at random and, conditionally on v , let $\theta = \theta_v$. Justify

$$\mathfrak{M}(\mathcal{L}_{\text{ada}}^n(\mathbb{F}^2), \|\cdot\|_2^2) \geq \inf_{\hat{\theta}} \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \mathbb{E}_{\theta_v} [\|\hat{\theta} - \theta_v\|_2^2].$$

Argue it is no loss of generality to assume that the choices for X_i are deterministic based on the Y_1^{i-1} . Thus, throughout we assume that $X_i = F_i(Y_1^{i-1}, u_1^i)$, where u_1^n is a fixed sequence, or, for simplicity, that X_i is a function of Y_1^{i-1} .

- (ii) Fix $\delta > 0$. Let $v \in \{-1, 1\}^d$, and for each such v , define $\theta_v = \delta v$. Also let P_v^n denote the joint distribution (over all adaptively chosen X_i) of the observed variables Y_1, \dots, Y_n , and define $P_{+j}^n = \frac{1}{2^{d-1}} \sum_{v: v_j=1} P_v^n$ and $P_{-j}^n = \frac{1}{2^{d-1}} \sum_{v: v_j=-1} P_v^n$, so that $P_{\pm j}^n$ denotes the distribution of the Y_i when $v \in \{-1, 1\}^d$ is chosen uniformly at random but conditioned on $v_j = \pm 1$. Then

$$\inf_{\hat{\theta}} \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \mathbb{E}_{\theta_v} [\|\hat{\theta} - \theta_v\|_2^2] \geq \frac{\delta^2}{2} \sum_{j=1}^d \left[1 - \|P_{+j}^n - P_{-j}^n\|_{\text{TV}} \right].$$

- (iii) We have

$$\frac{\delta^2}{2} \sum_{j=1}^d \left[1 - \|P_{+j}^n - P_{-j}^n\|_{\text{TV}} \right] \geq \frac{\delta^2 d}{2} \left[1 - \left(\frac{1}{d} \sum_{j=1}^d \|P_{+j}^n - P_{-j}^n\|_{\text{TV}}^2 \right)^{\frac{1}{2}} \right].$$

- (iv) Let $P_{+j}^{(i)}$ be the distribution of the random variable Y_i conditioned on $v_j = +1$ (with the other coordinates of v chosen uniformly at random), and let $P_{+j}^{(i)}(\cdot \mid y_1^{i-1}, x_i)$ denote the distribution of Y_i conditioned on $v_j = +1$, $Y_1^{i-1} = y_1^{i-1}$, and x_i . Justify

$$\begin{aligned} \|P_{+j}^n - P_{-j}^n\|_{\text{TV}}^2 &\leq \frac{1}{2} D_{\text{kl}}(P_{+j}^n \| P_{-j}^n) \\ &\leq \frac{1}{2} \sum_{i=1}^n \int D_{\text{kl}}(P_{+j}^{(i)}(\cdot \mid y_1^{i-1}, x_i) \| P_{-j}^{(i)}(\cdot \mid y_1^{i-1}, x_i)) dP_{+j}^{i-1}(y_1^{i-1}, x_i). \end{aligned}$$

- (v) Then we have

$$\sum_{j=1}^d D_{\text{kl}}(P_{+j}^{(i)}(\cdot \mid y_1^{i-1}, x_i) \| P_{-j}^{(i)}(\cdot \mid y_1^{i-1}, x_i)) \leq \frac{2\delta^2}{\sigma^2} \|x_i\|_2^2.$$

- (vi) We have

$$\sum_{j=1}^d \|P_{+j}^n - P_{-j}^n\|_{\text{TV}}^2 \leq \frac{\delta^2}{\sigma^2} \mathbb{E}[\|X\|_{\text{Fr}}^2],$$

where the final expectation is over V drawn uniformly in $\{-1, 1\}^d$ and all Y_i, X_i .

- (vii) Show how to choose δ appropriately to conclude the minimax bound (9.8.5).

Exercise 9.7: Suppose under the setting of Question 9.6 that we may no longer be adaptive, meaning that the matrix $X \in \mathbb{R}^{n \times d}$ must be chosen ahead of time (without seeing any data). Assuming $n \geq d$, is it possible to attain (within a constant factor) the risk (9.8.5)? If so, give an example construction, if not, explain why not.

JCD Comment: Put this in the next chapter

Exercise 9.8 (The curse of dimensionality in nonparametric regression): Consider the nonparametric regression problem in Section 10.1. Let \mathbb{B}^d be the unit ℓ_2 -ball in \mathbb{R}^d and consider the function class \mathcal{F} of 1-Lipschitz functions taking values in $[-1, 1]$ on \mathbb{B}^d , and consider the error $\|f - g\|_2^2 = \int_{\mathbb{B}^d} (f(x) - g(x))^2 dx$. (Here, 1-Lipschitz means $|f(x) - f(x')| \leq \|x - x'\|_2$ for any x, x' .) We show the minimax lower bound (10.1.4) for this function class using Fano's method. Fix $\delta \in [0, 1]$ to be chosen and let $\{x_j\}_{j=1}^M$ be the centers of a maximal 2δ -packing of \mathbb{B}^d , so that $M \geq (\frac{1}{2\delta})^d$ (by Lemma 5.1.10), and define the “bump” functions

$$g_j(x) = \delta [1 - \|x - x_j\|_2 / \delta]_+,$$

which all have disjoint support. Then for a vector $v \in \{\pm 1\}^M$, define

$$f_v(x) := \sum_{j=1}^M v_j g_j(x).$$

- (a) Show that $f_v \in \mathcal{F}$.
- (b) Show that $\int g_j(x)^2 dx = \frac{2 \cdot \text{SA}(d)}{d(d+1)(d+2)} \delta^{2+d}$, where $\text{SA}(d)$ denotes the surface area of \mathbb{B}^d .
- (c) Use the Gilbert-Varshamov bound (Lemma 9.2.3) to show there is a collection $\mathcal{V} \subset \{\pm 1\}^M$ of cardinality $\exp(M/8)$ with $\|f_v - f_{v'}\|_2^2 \geq c_d \delta^2$ for all $v \neq v' \in \mathcal{V}$, where c_d depends only on the dimension d .
- (d) Prove the minimax lower bound (10.1.4) for $\beta = 1$.

Exercise 9.9 (Optimal algorithms for memory access): In a modern CPU, memory is organized in a hierarchy, so that data upon which computations are being actively performed lies in a very small memory close to the logic units of the processor for which access is extraordinarily fast, while data not being actively used lies in slower memory slightly farther from the processor. (Modern processor memory is generally organized into the registers—a small number of 4- or 8-byte memory locations on the processor—and level 1, 2, (and sometimes 3 or more) cache, which contain small amounts of data and increasing access times, and RAM (random access memory).) Moving data—communicating—between levels of the memory hierarchy is both power intensive and very slow relative to computation on the data itself, so that in many algorithms the bulk of the time of the algorithm is in moving data from one place to another to be computed upon. Thus, developing very fast algorithms for numerical (and other) tasks on modern computers requires careful tracking of memory access and communication, and careful control of these quantities can often yield orders of magnitude speed improvements in execution. In this problem, you will prove a lower bound on the number of communication steps that a variety of numerical-type methods must perform, giving a concrete (attainable) inequality that allows one to certify optimality of *specific* algorithms.

In particular, we consider matrix multiplication, as it is a proxy for a class of cubic algorithms that are well behaved. Let $A, B \in \mathbb{R}^{n \times n}$ be matrices, and assume we wish to compute $C = AB$, via the simple algorithm that for all i, j sets

$$C_{ij} = \sum_{l=1}^n A_{il} B_{lj}.$$

Computationally, this forces us to repeatedly execute operations of the form

$$\text{Mem}(C_{ij}) = F(\text{Mem}(A_{il}), \text{Mem}(B_{lj}), \text{Mem}(C_{ij})),$$

where F is some function—that may depend on i, j, l —and $\text{Mem}(\cdot)$ indicates that we access the memory associated with the argument. (In our case, we have $C_{ij} = C_{ij} + A_{il} \cdot B_{lj}$.) We assume that executing F requires that $\text{Mem}(A_{il})$, $\text{Mem}(B_{lj})$, and $\text{Mem}(C_{ij})$ belong to fast memory, and that each are distinct (stored in a separate place in flow and fast memory). We assume that the order of the computations does *not* matter, so we may re-order them in any way. We call $\text{Mem}(A_{il})$ (respectively B or C) and *operand* in our computation. We let M denote the size of fast/local memory, and we would like to lower bound the number of times we must communicate an operand into or out of the fast local memory as a function of n , the matrix size, and M , the fast memory size, when all we may do is re-order the computation being executed. We let N_{Store} denote the number of times we write something from fast memory out to slow memory and let N_{Load} the number of times we load something from slow memory to fast memory. Let N be the total number of operations we execute (for simple matrix multiplication, we have $N = n^3$, though with sparse matrices, this can be smaller).

We analyze the procedure by breaking the computation into a number of segments, where each segment contains precisely M load or store (communication-causing) instructions.

- (a) Let N_{seg} be an upper bound on the number of evaluations with the function $F(\cdot)$ in any given segment (you will upper bound this in a later part of the problem). Justify that

$$N_{\text{Store}} + N_{\text{Load}} \geq M \lfloor N/N_{\text{seg}} \rfloor.$$

- (b) Within a segment, all operands involved must be in fast memory at least once to be computed with. Assume that memory locations $\text{Mem}(A_{il})$, $\text{Mem}(B_{lj})$, and $\text{Mem}(C_{ij})$ do not overlap. For any operand involved in a memory operation in one of the segments, the operand (1) was already in fast memory at the beginning of the segment, (2) was read from slow memory, (3) is still in fast memory at the end of the segment, or (4) is written to slow memory at the end of the segment. (There are also operands potentially created during execution that are simply discarded; we do not bound those.) Justify the following: within a segment, for each type of operand $\text{Mem}(A_{ij})$, $\text{Mem}(B_{ij})$, or $\text{Mem}(C_{ij})$, there are at most $c \cdot M$ such operands (i.e. there are at most cM operands of type $\text{Mem}(A_{ij})$, independent of the others, and so on), where c is a numerical constant. What value of c can you attain?
- (c) Using the result of question 7.2, argue that $N_{\text{seg}} \leq c' \sqrt{M^3}$ for a numerical constant c' . What value of c' do you get?
- (d) Using the result of part (c), argue that the number of loads and stores satisfies

$$N_{\text{Store}} + N_{\text{Load}} \geq c'' \frac{N}{\sqrt{M}} - M$$

for a numerical constant c'' . What is your constant?

Exercise 9.10 (Tilting and information bounds): Let P_0 be a distribution on $X \in \mathcal{X}$, and assume w.l.o.g. that P_0 has a density p_0 . For a bounded function $g : \mathcal{X} \rightarrow \mathbb{R}^k$ with $\mathbb{E}_{P_0}[g(X)] = 0$, define the *tilted density*

$$p_{t,g}(x) := (1 + \langle t, g(x) \rangle) p_0(x), \quad (9.8.6)$$

for $t \in \mathbb{R}^k$, with induced distribution $P_{t,g}$.

(a) Show that if t is near enough to 0, then $p_{t,g}$ is indeed a valid density.

(b) Show that for $t \in \mathbb{R}^k$ near 0,

$$D_{\text{kl}}(P_{t,g} \| P_0) = \frac{1}{2} t^\top \text{Cov}_0(g(X)) t + O(\|t\|^3) \quad \text{and} \quad D_{\text{kl}}(P_0 \| P_{t,g}) = \frac{1}{2} t^\top \text{Cov}_0(g(X)) t + O(\|t\|^3).$$

(c) Let $\mathcal{V} \subset \mathbb{R}^k$ be a bounded set and $V \sim \text{Uniform}(\mathcal{V})$. Let $\delta \geq 0$. Suppose that conditional on $V = v$, we draw $X_i \stackrel{\text{iid}}{\sim} P_{\delta v, g}$. Show the mutual information bound

$$I(X_1^n; V) \leq \frac{n\delta^2}{2} \text{tr} \left(\text{Cov}_0(g(X)) \mathbb{E}[VV^\top] \right) + O(n\delta^3).$$

Exercise 9.11 (General tilting): Replace the tilting (9.8.6) with

$$p_{t,g}(x) = \frac{1}{C_t} \phi(\langle t, g(x) \rangle) p_0(x), \quad C_t := \int \phi(\langle t, g(x) \rangle) p_0(x) d\mu(x)$$

where $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$ is a bounded nonnegative function, differentiable at 0 with $\phi(0) = \phi'(0) = 1$, with bounded second derivatives. (For example, $\phi(t) = \frac{2}{1+e^{-2t}}$.) Show that if $\mathbb{E}_0[\|g(X)\|^2] < \infty$ then the following analogues of the results (a), (b), and (c) from Exercise 9.10 hold:

(a) The normalization $C_t := \mathbb{E}_0[\phi(\langle t, g(X) \rangle)]$ satisfies $C_t = 1 + \frac{\phi''(0)}{2} t^\top \text{Cov}_0(g(X)) t + o(\|t\|^2)$.

(b) For t near 0,

$$D_{\text{kl}}(P_{t,g} \| P_0) = \frac{1}{2} t^\top \text{Cov}_0(g(X)) t + o(\|t\|^2).$$

(c) Let $\mathcal{V} \subset \mathbb{R}^k$ be a bounded set and $V \sim \text{Uniform}(\mathcal{V})$. Let $\delta \geq 0$. Suppose that conditional on $V = v$, we draw $X_i \stackrel{\text{iid}}{\sim} P_{\delta v, g}$. Then

$$I(X_1^n; V) \leq \frac{n\delta^2}{2} \text{tr} \left(\text{Cov}_0(g(X)) \mathbb{E}[VV^\top] \right) + n \cdot o(\delta^2).$$

Exercise 9.12: Let \mathcal{P} be a collection of distributions on a set \mathcal{X} , and let $\theta : \mathcal{P} \rightarrow \mathbb{R}^d$ be a parameter of the distribution. Let \mathcal{G} a collection of mean-zero functions mapping $\mathcal{X} \rightarrow \mathbb{R}^k$ with $\mathbb{E}_0[g(X)] = 0$ for each $g \in \mathcal{G}$. For $t \in \mathbb{R}^k$ define the distribution $P_{t,g}$ by the tilting (9.8.6), and assume $P_{t,g} \in \mathcal{P}$ for small enough t . We say that θ is *differentiable at P_0 relative to \mathcal{G}* with derivative $\dot{\theta}_{P_0} : \mathcal{X} \rightarrow \mathbb{R}^d$ if

$$\theta(P_{t,g}) = \theta(P_0) + \mathbb{E}_{P_0}[\dot{\theta}_{P_0}(X) \langle g(X), t \rangle] + o(\|t\|)$$

as $t \rightarrow 0$ for each $g \in \mathcal{G}$, where $\mathbb{E}_{P_0}[\dot{\theta}_{P_0}(X)] = 0$.

- (a) Show that the mean is differentiable relative to \mathcal{G} for any collection of mean-zero bounded functions, and give the derivative $\dot{\theta}_{P_0}$.

For the remainder of the problem, assume that for any \mathcal{G} consisting of bounded functions, $\theta : \mathcal{P} \rightarrow \mathbb{R}^d$ is differentiable relative to \mathcal{G} with bounded derivative $\dot{\theta}_{P_0}$. Let $\Sigma = \text{Cov}_0(\dot{\theta}_{P_0})$ and for $A \succ 0$ define the Mahalanobis norm $\|v\|_A^2 = v^\top A v$.

- (b) For $\delta \geq 0$, show how to construct a collection of distributions $\{Q_v\}$ indexed by $v \in \mathcal{V} = \{-1, 1\}^d$ such that as $\delta \rightarrow 0$, for any pair v, v'

$$\|\theta(Q_v) - \theta(Q_{v'})\|_A = \delta \|v - v'\|_A + o(\delta).$$

- (c) Show that (for large enough dimension d) there exist numerical constants $c_0, c_1 > 0$ and a collection of distributions $\{Q_v\} \subset \mathcal{P}$, indexed by $v \in \mathcal{V} = \{-1, 1\}^d$ such that if we choose $V \sim \text{Uniform}(\mathcal{V})$ then draw $X_1^n \stackrel{\text{iid}}{\sim} Q_v$ conditional on $V = v$,

$$\mathbb{P}\left(\|\hat{\theta}(X_1^n) - \theta(Q_V)\|_{\Sigma^{-1}} \geq c_0 \sqrt{d/n}\right) \geq c_1.$$

Exercise 9.13 (A minimax lower bound for M-estimation): Let ℓ be a convex loss and $L(\theta) = \mathbb{E}_{P_0}[\ell(\theta, Z)]$ satisfy all the conditions of Assumption A.5.1 in Chapter 5.3, where P_0 is a distribution on $Z \in \mathcal{Z}$ that we assume w.l.o.g. has a density p_0 . Define the tilted densities p_t as in Exercise 9.10, Eq. (9.8.6). Let $\theta(P) = \text{argmin}_\theta \mathbb{E}_P[\ell(\theta, Z)]$ be the population minimizer and define the Hessian $H = \nabla^2 L(\theta_0)$ and covariance $\Sigma := \text{Cov}_{P_0}(\nabla \ell(\theta_0, Z)) \succ 0$. Previewing Chapter 12.3.1, the implicit function theorem implies that

$$\theta(P_{t,g}) = \theta(P_0) + H^{-1} \mathbb{E}_{P_0}[\nabla \ell(\theta_0, Z) g(Z)^\top] t + o(\|t\|)$$

for $t \in \mathbb{R}^k$ near 0. Show that by choosing an appropriate tilting function $g : \mathcal{Z} \rightarrow \mathbb{R}^d$, there exist numerical constants $c_0, c_1, c_2 > 0$ such for large enough n , the choice $\delta_n := c_0 \sqrt{\frac{d}{n}}$ yields

$$\inf_{\hat{\theta}} \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P_{\delta_n v, g} \left((\hat{\theta}(Z_1^n) - \theta(P_t))^\top (H^{-1} \Sigma H^{-1})^{-1} (\hat{\theta}(Z_1^n) - \theta(P_t)) \geq c_1 \frac{d}{n} \right) \geq c_2.$$

How does this lower bound compare to the guarantee in Corollary 5.3.10?

JCD Comment: A few additional question ideas:

1. Use the global Fano method technique to give lower bounds for density estimation
2. Curse of dimensionality in high-dimensional regression? The idea would be to take disjoint δ -balls $B_j \subset \mathbb{B}^d$, where $\mathbb{B}^d = \{x \mid \|x\| \leq 1\}$ is the unit ball, with centers x_j , where j runs from 1 to $(1/\delta)^d$, then define the bump function $g_j(x) = \delta [1 - \|x - x_j\| / \delta]_+$. Then set $f_v(x) = \sum_j v_j g_j(x)$, which is 1-Lipschitz for the norm $\|\cdot\|$. Then the separation is δ , while the log cardinality is $2^{\delta^{-d}}$, giving $\delta^2(1 - n\delta^{2+d})$ as the lower bound. Take $\delta = n^{-1/(2+d)}$.

Chapter 10

Beyond local minimax techniques

10.1 Nonparametric regression: minimax upper and lower bounds

To show further applications of the minimax optimality ideas we have developed, we consider one of the two the most classical non-parametric (meaning that the number of parameters can grow with the sample size n) problems: estimating a regression function on a subset of the real line (the most classical problem being estimation of a density). In non-parametric regression, we assume there is an unknown function $f : \mathbb{R} \rightarrow \mathbb{R}$, where f belongs to a pre-determined class of functions \mathcal{F} ; usually this class is parameterized by some type of smoothness guarantee. To make our problems concrete, we will assume that the unknown function f is L -Lipschitz and defined on $[0, 1]$. Let \mathcal{F} denote this class.

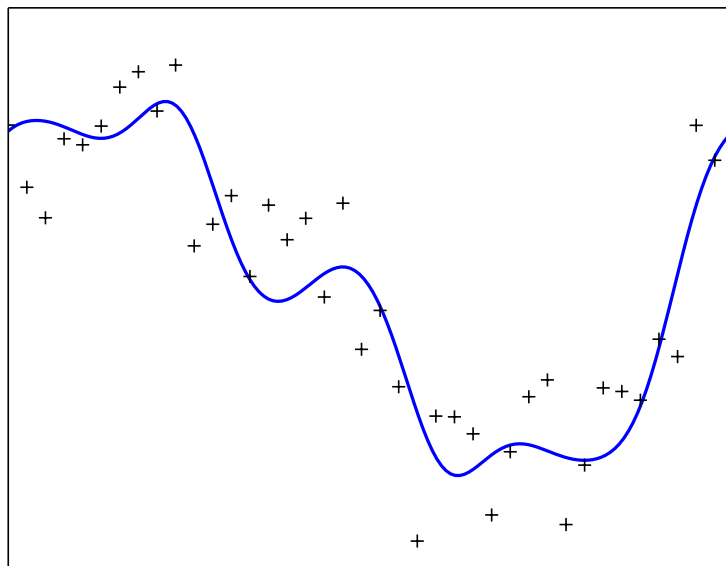


Figure 10.1. Observations in a non-parametric regression problem, with function f plotted. (Here $f(x) = \sin(2x + \cos^2(3x))$.)

In the standard non-parametric regression problem, we obtain observations of the form

$$Y_i = f(X_i) + \varepsilon_i \tag{10.1.1}$$

where ε_i are independent, mean zero conditional on X_i , and $\mathbb{E}[\varepsilon_i^2] \leq \sigma^2$. See Figure 10.1 for an example. We also assume that we fix the locations of the X_i as $X_i = i/n \in [0, 1]$, that is, the X_i are evenly spaced in $[0, 1]$. Given n observations Y_i , we ask two questions: (1) how can we estimate f ? and (2) what are the optimal rates at which it is possible to estimate f ?

10.1.1 Kernel estimates of the function

A natural strategy is to place small “bumps” around the observed points, and estimate f in a neighborhood of a point x by weighted averages of the Y values for other points near x . We now formalize a strategy for doing this. Suppose we have a kernel function $K : \mathbb{R} \rightarrow \mathbb{R}_+$, which is continuous, not identically zero, has support $\text{supp } K = [-1, 1]$, and satisfies the technical condition

$$\lambda_0 \sup_x K(x) \leq \inf_{|x| \leq 1/2} K(x), \quad (10.1.2)$$

where $\lambda_0 > 0$ (this says the kernel has some width to it). A natural example is the “tent” function given by $K_{\text{tent}}(x) = [1 - |x|]_+$, which satisfies inequality (10.1.2) with $\lambda_0 = 1/2$. See Fig. 10.2 for two examples, one the tent function and the other the function

$$K(x) = \mathbf{1}\{|x| < 1\} \exp\left(-\frac{1}{(x-1)^2}\right) \exp\left(-\frac{1}{(x+1)^2}\right),$$

which is infinitely differentiable and supported on $[-1, 1]$.

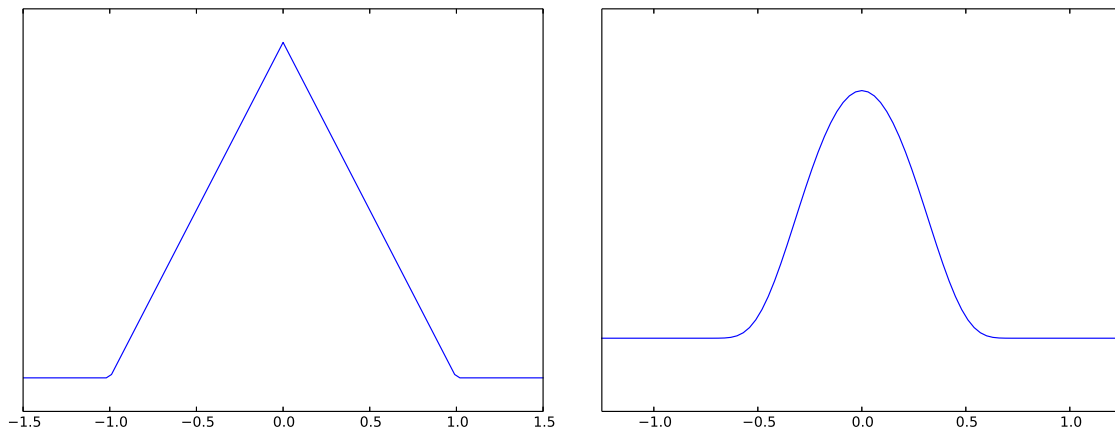


Figure 10.2: Left: “tent” kernel. Right: infinitely differentiable compactly supported kernel.

Now we consider a natural estimator of the function f based on observations (10.1.2) known as the Nadaraya-Watson estimator. Fix a bandwidth h , which we will see later smooths the estimated functions f . For all x , define weights

$$W_{ni}(x) := \frac{K\left(\frac{X_i - x}{h}\right)}{\sum_{j=1}^n K\left(\frac{X_j - x}{h}\right)}$$

and define the estimated function

$$\hat{f}_n(x) := \sum_{i=1}^n Y_i W_{ni}(x).$$

The intuition here is that we have a locally weighted regression function, where points X_i in the neighborhood of x are given higher weight than further points. Using this function \hat{f}_n as our estimator, it is possible to provide a guarantee on the bias and variance of the estimated function at each point $x \in [0, 1]$.

Proposition 10.1.1. *Let the observation model (10.1.1) hold and assume condition (10.1.2). In addition assume the bandwidth is suitably large that $h \geq 2/n$ and that the X_i are evenly spaced on $[0, 1]$. Then for any $x \in [0, 1]$, we have*

$$|\mathbb{E}[\hat{f}_n(x)] - f(x)| \leq Lh \quad \text{and} \quad \text{Var}(\hat{f}_n(x)) \leq \frac{2\sigma^2}{\lambda_0 n h}.$$

Proof To bound the bias, we note that (conditioning implicitly on X_i)

$$\mathbb{E}[\hat{f}_n(x)] = \sum_{i=1}^n \mathbb{E}[Y_i W_{ni}(x)] = \sum_{i=1}^n \mathbb{E}[f(X_i) W_{ni}(x) + \varepsilon_i W_{ni}(x)] = \sum_{i=1}^n f(X_i) W_{ni}(x).$$

Thus we have that the bias is bounded as

$$\begin{aligned} \left| \mathbb{E}[\hat{f}_n(x)] - f(x) \right| &\leq \sum_{i=1}^n |f(X_i) - f(x)| W_{ni}(x) \\ &\leq \sum_{i: |X_i - x| \leq h} |f(X_i) - f(x)| W_{ni}(x) \leq Lh \sum_{i=1}^n W_{ni}(x) = Lh. \end{aligned}$$

To bound the variance, we claim that

$$W_{ni}(x) \leq \min \left\{ \frac{2}{\lambda_0 n h}, 1 \right\}. \quad (10.1.3)$$

Indeed, we have that

$$W_{ni}(x) = \frac{K\left(\frac{X_i - x}{h}\right)}{\sum_{j=1}^n K\left(\frac{X_j - x}{h}\right)} = \frac{K\left(\frac{X_i - x}{h}\right)}{\sum_{j: |X_j - x| \leq h/2} K\left(\frac{X_j - x}{h}\right)} \leq \frac{K\left(\frac{X_i - x}{h}\right)}{\lambda_0 \sup_x K(x) |\{j : |X_j - x| \leq h/2\}|},$$

and because there are at least $nh/2$ indices satisfying $|X_j - x| \leq h$, we obtain the claim (10.1.3). Using the claim, we have

$$\begin{aligned} \text{Var}(\hat{f}_n(x)) &= \mathbb{E} \left[\left(\sum_{i=1}^n (Y_i - f(X_i)) W_{ni}(x) \right)^2 \right] = \mathbb{E} \left[\left(\sum_{i=1}^n \varepsilon_i W_{ni}(x) \right)^2 \right] \\ &= \sum_{i=1}^n W_{ni}(x)^2 \mathbb{E}[\varepsilon_i^2] \leq \sum_{i=1}^n \sigma^2 W_{ni}(x)^2. \end{aligned}$$

Noting that $W_{ni}(x) \leq 2/\lambda_0 n h$ and $\sum_{i=1}^n W_{ni}(x) = 1$, we have

$$\sum_{i=1}^n \sigma^2 W_{ni}(x)^2 \leq \sigma^2 \max_i W_{ni}(x) \underbrace{\sum_{i=1}^n W_{ni}(x)}_{=1} \leq \sigma^2 \frac{2}{\lambda_0 n h},$$

completing the proof. \square

With the proposition in place, we can then provide a theorem bounding the worst case pointwise mean squared error for estimation of a function $f \in \mathcal{F}$.

Theorem 10.1.2. *Under the conditions of Proposition 10.1.1, choose $h = (\sigma^2/L^2\lambda_0)^{1/3}n^{-1/3}$. Then there exists a universal (numerical) constant $C < \infty$ such that for any $f \in \mathcal{F}$,*

$$\sup_{x \in [0,1]} \mathbb{E}[(\hat{f}_n(x) - f(x))^2] \leq C \left(\frac{L\sigma^2}{\lambda_0} \right)^{2/3} n^{-\frac{2}{3}}.$$

Proof Using Proposition 10.1.1, we have for any $x \in [0, 1]$ that

$$\mathbb{E}[(\hat{f}_n(x) - f(x))^2] = \left(\mathbb{E}[\hat{f}_n(x)] - f(x) \right)^2 + \mathbb{E}[(\hat{f}_n(x) - \mathbb{E}[\hat{f}_n(x)])^2] \leq \frac{2\sigma^2}{\lambda_0 n h} + L^2 h^2.$$

Choosing h to balance the above bias/variance tradeoff, we obtain the theorem. \square

By integrating the result in Theorem 10.1.2 over the interval $[0, 1]$, we immediately obtain the following corollary.

Corollary 10.1.3. *Under the conditions of Theorem 10.1.2, if we use the tent kernel K_{tent} , we have*

$$\sup_{f \in \mathcal{F}} \mathbb{E}_f[\|\hat{f}_n - f\|_2^2] \leq C \left(\frac{L\sigma^2}{n} \right)^{2/3},$$

where C is a universal constant.

In Proposition 10.1.1, it is possible to show that a more clever choice of kernels—ones that are not always positive—can attain bias $\mathbb{E}[\hat{f}_n(x)] - f(x) = \mathcal{O}(h^\beta)$ if f has Lipschitz $(\beta - 1)$ th derivative. In this case, we immediately obtain that the rate can be improved to

$$\sup_x \mathbb{E}[(\hat{f}_n(x) - f(x))^2] \leq C n^{-\frac{2\beta}{2\beta+1}},$$

and every additional degree of smoothness gives a corresponding improvement in convergence rate. We also remark that rates of this form, which are much larger than n^{-1} , are characteristic of non-parametric problems; essentially, we must adaptively choose a dimension that balances the sample size, so that rates of $1/n$ are difficult or impossible to achieve.

10.1.2 Minimax lower bounds on estimation with Assouad's method

Now we can ask whether the results we have given are in fact sharp; do there exist estimators attaining a faster rate of convergence than our kernel-based (locally weighted) estimator? Using Assouad's method, we show that, in fact, these results are all tight. In particular, we prove the following result on minimax estimation of a regression function $f \in \mathcal{F}$, where \mathcal{F} consists of 1-Lipschitz functions defined on $[0, 1]$, in the $\|\cdot\|_2^2$ error, that is, $\|f - g\|_2^2 = \int_0^1 (f(t) - g(t))^2 dt$.

Theorem 10.1.4. *Let the observation points X_i be spaced evenly on $[0, 1]$, and assume the observation model (10.1.1). Then there exists a universal constant $c > 0$ such that*

$$\mathfrak{M}_n(\mathcal{F}, \|\cdot\|_2^2) := \inf_{\hat{f}_n} \sup_{f \in \mathcal{F}} \mathbb{E}_f \left[\|\hat{f}_n - f\|_2^2 \right] \geq c \left(\frac{\sigma^2}{n} \right)^{\frac{2}{3}}.$$

Deferring the proof of the theorem temporarily, we make a few remarks. It is in fact possible to show—using a completely identical technique—that if \mathcal{F}_β denotes the class of functions with $\beta - 1$ derivatives, where the $(\beta - 1)$ th derivative is Lipschitz, then

$$\mathfrak{M}_n(\mathcal{F}_\beta, \|\cdot\|_2^2) \geq c \left(\frac{\sigma^2}{n} \right)^{\frac{2\beta}{2\beta+1}}.$$

So for any smoothness class, we can never achieve the parametric σ^2/n rate, but we can come arbitrarily close. As another remark, which we do not prove, in dimensions $d \geq 1$, the minimax rate for estimation of functions f with Lipschitz $(\beta - 1)$ th derivative scales as

$$\mathfrak{M}_n(\mathcal{F}_\beta, \|\cdot\|_2^2) \geq c \left(\frac{\sigma^2}{n} \right)^{\frac{2\beta}{2\beta+d}}. \quad (10.1.4)$$

This result can, similarly, be proved using a variant of Assouad’s method or a local Fano method; see, for example, Györfi et al. [108, Chapter 3]. Exercise 9.8 works through a particular case of this lower bound. This is a striking example of the curse of dimensionality: the penalty for increasing dimension results in worse rates of convergence. For example, suppose that $\beta = 1$. In 1 dimension, we require $n \geq 90 \approx (.05)^{-3/2}$ observations to achieve accuracy .05 in estimation of f , while we require $n \geq 8000 = (.05)^{-(2+d)/2}$ even when the dimension $d = 4$, and $n \geq 64 \cdot 10^6$ observations even in 10 dimensions, which is a relatively small problem. That is, the problem is made exponentially more difficult by dimension increases.

We now prove Theorem 10.1.4. To establish the result, we show how to construct a family of problems—indexed by binary vectors $v \in \{-1, 1\}^k$ —so that our estimation problem satisfies the separation (9.5.1), then we show that the information based on observing noisy versions of the functions we have defined is small. Choosing k to make our resulting lower bound as high as possible completes the argument.

Construction of a separated family of functions To construct our separation in Hamming metric, as required by Eq. (9.5.1), fix some $k \in \mathbb{N}$; we will choose k later. This approach is somewhat different from our standard approach of using a fixed dimensionality and scaling the separation directly; in non-parametric problems, we scale the “dimension” itself to adjust the difficulty of the estimation problem. Define the function $g(x) = [1/2 - |x - 1/2|]_+$, so that g is 1-Lipschitz and is 0 outside of the interval $[0, 1]$. Then for any $v \in \{-1, 1\}^k$, define the “bump” functions

$$g_j(x) := \frac{1}{k} g \left(k \left(x - \frac{j-1}{k} \right) \right) \quad \text{and} \quad f_v(x) := \sum_{j=1}^k v_j g_j(x),$$

which we see is 1-Lipschitz. Now, consider any function $f : [0, 1] \rightarrow \mathbb{R}$, and let E_j be shorthand for the intervals $E_j = [(j-1)/k, j/k]$ for $j = 1, \dots, k$. We must find a mapping identifying a function f with points in the hypercube $\{-1, 1\}^k$. To that end, we may define a vector $\hat{v}(f) \in \{-1, 1\}^k$ by

$$\hat{v}_j(f) = \operatorname{argmin}_{s \in \{-1, 1\}} \int_{E_j} (f(t) - sg_j(t))^2 dt.$$

We claim that for any function f ,

$$\left(\int_{E_j} (f(t) - f_v(t))^2 dt \right)^{\frac{1}{2}} \geq \mathbf{1} \{ \hat{v}_j(f) \neq v_j \} \left(\int_{E_j} f_v(t)^2 dt \right)^{\frac{1}{2}}. \quad (10.1.5)$$

Indeed, on the set E_j , we have $v_j g_j(t) = f_v(t)$, and thus $\int_{E_j} g_j(t)^2 dt = \int_{E_j} f_v(t)^2 dt$. Then by the triangle inequality, we have

$$\begin{aligned} 2 \cdot \mathbf{1}\{\widehat{v}_j(f) \neq v_j\} \left(\int_{E_j} g_j(t)^2 dt \right)^{\frac{1}{2}} &= \left(\int_{E_j} ((\widehat{v}_j(f) - v_j) g_j(t))^2 dt \right)^{\frac{1}{2}} \\ &\leq \left(\int_{E_j} (f(t) - v_j g_j(t))^2 dt \right)^{\frac{1}{2}} + \left(\int_{E_j} (f(t) - \widehat{v}_j(f) g_j(t))^2 dt \right)^{\frac{1}{2}} \\ &\leq 2 \left(\int_{E_j} (f(t) - f_v(t))^2 dt \right)^{\frac{1}{2}}, \end{aligned}$$

by definition of the sign $\widehat{v}_j(f)$.

With the definition of \widehat{v} and inequality (10.1.5), we see that for any vector $v \in \{-1, 1\}^k$, we have

$$\|f - f_v\|_2^2 = \sum_{j=1}^k \int_{E_j} (f(t) - f_v(t))^2 dt \geq \sum_{j=1}^k \mathbf{1}\{\widehat{v}_j(f) \neq v_j\} \int_{E_j} f_v(t)^2 dt.$$

In particular, we know that

$$\int_{E_j} f_v(t)^2 dt = \frac{1}{k^2} \int_0^{1/k} g(kt)^2 dt = \frac{1}{k^3} \int_0^1 g(u)^2 du \geq \frac{c}{k^3},$$

where c is a numerical constant. In particular, we have the desired separation

$$\|f - f_v\|_2^2 \geq \frac{c}{k^3} \sum_{j=1}^k \mathbf{1}\{\widehat{v}_j(f) \neq v_j\}. \quad (10.1.6)$$

Bounding the binary testing error Let P_v^n denote the distribution of the n observations $Y_i = f_v(X_i) + \varepsilon_i$ when f_v is the true regression function. Then inequality (10.1.6) implies via Assouad's lemma that

$$\mathfrak{M}_n(\mathcal{F}, \|\cdot\|_2^2) \geq \frac{c}{k^3} \sum_{j=1}^k \left[1 - \|P_{+j}^n - P_{-j}^n\|_{\text{TV}} \right]. \quad (10.1.7)$$

Now, we use convexity and Pinsker's inequality to note that

$$\|P_{+j}^n - P_{-j}^n\|_{\text{TV}}^2 \leq \max_v \|P_{v,+j}^n - P_{v,-j}^n\|_{\text{TV}}^2 \leq \max_v \frac{1}{2} D_{\text{kl}}(P_{v,+j}^n \| P_{v,-j}^n).$$

For any two functions f_v and $f_{v'}$, we have that the observations Y_i are independent and normal with means $f_v(X_i)$ or $f_{v'}(X_i)$, respectively. Thus

$$\begin{aligned} D_{\text{kl}}(P_v^n \| P_{v'}^n) &= \sum_{i=1}^n D_{\text{kl}}(\mathcal{N}(f_v(X_i), \sigma^2) \| \mathcal{N}(f_{v'}(X_i), \sigma^2)) \\ &= \sum_{i=1}^n \frac{1}{2\sigma^2} (f_v(X_i) - f_{v'}(X_i))^2. \end{aligned} \quad (10.1.8)$$

Now we must show that the expression (10.1.8) scales more slowly than n , which we will see must be the case as whenever $d_{\text{ham}}(v, v') \leq 1$. Intuitively, most of the observations have the same distribution by our construction of the f_v as bump functions; let us make this rigorous.

We may assume without loss of generality that $v_j = v'_j$ for $j > 1$. As the $X_i = i/n$, we thus have that only X_i for i near 1 can have non-zero values in the tensorization (10.1.8). In particular,

$$f_v(i/n) = f_{v'}(i/n) \quad \text{for all } i \text{ s.t. } \frac{i}{n} \geq \frac{2}{k}, \quad \text{i.e. } i \geq \frac{2n}{k}.$$

Rewriting expression (10.1.8), then, and noting that $f_v(x) \in [-1/k, 1/k]$ for all x by construction, we have

$$\sum_{i=1}^n \frac{1}{2\sigma^2} (f_v(X_i) - f_{v'}(X_i))^2 \leq \sum_{i=1}^{2n/k} \frac{1}{2\sigma^2} (f_v(X_i) - f_{v'}(X_i))^2 \leq \frac{1}{2\sigma^2} \frac{2n}{k} \frac{1}{k^2} = \frac{n}{k^3\sigma^2}.$$

Combining this with inequality (10.1.8) and the minimax bound (10.1.7), we obtain

$$\|P_{+j}^n - P_{-j}^n\|_{\text{TV}} \leq \sqrt{\frac{n}{2k^3\sigma^2}},$$

so

$$\mathfrak{M}_n(\mathcal{F}, \|\cdot\|_2^2) \geq \frac{c}{k^3} \sum_{j=1}^k \left[1 - \sqrt{\frac{n}{2k^3\sigma^2}} \right].$$

Choosing k for optimal tradeoffs Now we simply choose k ; in particular, setting

$$k = \left\lceil \left(\frac{n}{2\sigma^2} \right)^{1/3} \right\rceil \quad \text{then} \quad 1 - \sqrt{\frac{n}{2k^3\sigma^2}} \geq 1 - \sqrt{1/4} = \frac{1}{2},$$

and we arrive at

$$\mathfrak{M}_n(\mathcal{F}, \|\cdot\|_2^2) \geq \frac{c}{k^3} \sum_{j=1}^k \frac{1}{2} = \frac{c}{2k^2} \geq c' \left(\frac{\sigma^2}{n} \right)^{2/3},$$

where $c' > 0$ is a universal constant. Theorem 10.1.4 follows.

10.2 Global Fano Method

In this section, we extend the techniques of Section 9.4 on Fano's method (the local Fano method) to a more global construction. In particular, we show that, rather than constructing a local packing, choosing a scaling $\delta > 0$, and then optimizing over this δ , it is actually, in many cases, possible to prove lower bounds on minimax error directly using packing and covering numbers (metric entropy and packing entropy).

10.2.1 A mutual information bound based on metric entropy

To begin, we recall the classical Fano inequality in Corollary 9.4.2, which says that for any Markov chain $V \rightarrow X \rightarrow \hat{V}$, where V is uniform on the finite set \mathcal{V} , we have

$$\mathbb{P}(\hat{V} \neq V) \geq 1 - \frac{I(V; X) + \log 2}{\log(|\mathcal{V}|)}.$$

Thus, there are two ingredients in proving lower bounds on the error in a hypothesis test: upper bounding the mutual information and lower bounding the size $|\mathcal{V}|$. The key in the global Fano method is an upper bound on the former (the information $I(V; X)$) using covering numbers.

Before stating our result, we require a bit of notation. First, we assume that V is drawn from a distribution μ , and conditional on $V = v$, assume the sample $X \sim P_v$. Then a standard calculation (or simply the definition of mutual information; recall equation (9.4.4)) gives that

$$I(V; X) = \int D_{\text{kl}}(P_v \| \bar{P}) d\mu(v), \quad \text{where} \quad \bar{P} = \int P_v d\mu(v).$$

Now, we show how to connect this mutual information quantity to a covering number of a set of distributions.

Assume that for all v , we have $P_v \in \mathcal{P}$, where \mathcal{P} is a collection of distributions. In analogy with Definition 5.1, we say that the collection of distributions $\{Q_i\}_{i=1}^N$ form an ϵ -cover of \mathcal{P} in KL-divergence if for all $P \in \mathcal{P}$, there exists some i such that $D_{\text{kl}}(P \| Q_i) \leq \epsilon^2$. With this, we may define the KL-covering number of the set \mathcal{P} as

$$N_{\text{kl}}(\epsilon, \mathcal{P}) := \inf \left\{ N \in \mathbb{N} \mid \exists Q_i, i = 1, \dots, N, \sup_{P \in \mathcal{P}} \min_i D_{\text{kl}}(P \| Q_i) \leq \epsilon^2 \right\}, \quad (10.2.1)$$

where $N_{\text{kl}}(\epsilon, \mathcal{P}) = +\infty$ if no such cover exists. With definition (10.2.1) in place, we have the following proposition.

Proposition 10.2.1. *Under conditions of the preceding paragraphs, we have*

$$I(V; X) \leq \inf_{\epsilon > 0} \left\{ \epsilon^2 + \log N_{\text{kl}}(\epsilon, \mathcal{P}) \right\}. \quad (10.2.2)$$

Proof First, we claim that

$$\int D_{\text{kl}}(P_v \| \bar{P}) d\mu(v) \leq \int D_{\text{kl}}(P_v \| Q) d\mu(v) \quad (10.2.3)$$

for any distribution Q . Indeed, we have

$$\begin{aligned} \int D_{\text{kl}}(P_v \| \bar{P}) d\mu(v) &= \int_{\mathcal{V}} \int_{\mathcal{X}} dP_v \log \frac{dP_v}{d\bar{P}} d\mu(v) = \int_{\mathcal{V}} \int_{\mathcal{X}} dP_v \left[\log \frac{dP_v}{Q} + \log \frac{dQ}{d\bar{P}} \right] d\mu(v) \\ &= \int_{\mathcal{V}} D_{\text{kl}}(P_v \| Q) d\mu(v) + \underbrace{\int_{\mathcal{X}} \int_{\mathcal{V}} d\mu(v) dP_v \log \frac{dQ}{d\bar{P}}}_{=d\bar{P}} \\ &= \int D_{\text{kl}}(P_v \| Q) d\mu(v) - D_{\text{kl}}(\bar{P} \| Q) \leq \int D_{\text{kl}}(P_v \| Q) d\mu(v), \end{aligned}$$

so that inequality (10.2.3) holds. By carefully choosing the distribution Q in the upper bound (10.2.3), we obtain the proposition.

Now, assume that the distributions $Q_i, i = 1, \dots, N$ form an ϵ^2 -cover of the family \mathcal{P} , meaning that

$$\min_{i \in [N]} D_{\text{kl}}(P \| Q_i) \leq \epsilon^2 \quad \text{for all } P \in \mathcal{P}.$$

Let p_v and q_i denote the densities of P_v and Q_i with respect to some fixed base measure on \mathcal{X} (the choice of base measure does not matter). Then defining the distribution $Q = (1/N) \sum_{i=1}^N Q_i$, we obtain for any v that in expectation over $X \sim P_v$,

$$\begin{aligned} D_{\text{kl}}(P_v \| Q) &= \mathbb{E}_{P_v} \left[\log \frac{p_v(X)}{q(X)} \right] = \mathbb{E}_{P_v} \left[\log \frac{p_v(X)}{N^{-1} \sum_{i=1}^N q_i(X)} \right] \\ &= \log N + \mathbb{E}_{P_v} \left[\log \frac{p_v(X)}{\sum_{i=1}^N q_i(X)} \right] \leq \log N + \mathbb{E}_{P_v} \left[\log \frac{p_v(X)}{\max_i q_i(X)} \right] \\ &\leq \log N + \min_i \mathbb{E}_{P_v} \left[\log \frac{p_v(X)}{q_i(X)} \right] = \log N + \min_i D_{\text{kl}}(P_v \| Q_i). \end{aligned}$$

By our assumption that the Q_i form a cover, this gives the desired result, as $\epsilon \geq 0$ was arbitrary, as was our choice of the cover. \square

By a completely parallel proof, we also immediately obtain the following corollary.

Corollary 10.2.2. *Assume that X_1, \dots, X_n are drawn i.i.d. from P_v conditional on $V = v$. Let $N_{\text{kl}}(\epsilon, \mathcal{P})$ denote the KL-covering number of a collection \mathcal{P} containing the distributions (over a single observation) P_v for all $v \in \mathcal{V}$. Then*

$$I(V; X_1, \dots, X_n) \leq \inf_{\epsilon \geq 0} \{n\epsilon^2 + \log N_{\text{kl}}(\epsilon, \mathcal{P})\}.$$

With Corollary 10.2.2 and Proposition 10.2.1 in place, we thus see that the global covering numbers in KL-divergence govern the behavior of information.

We remark in passing that the quantity (10.2.2), and its i.i.d. analogue in Corollary 10.2.2, is known as the *index of resolvability*, and it controls estimation rates and redundancy of coding schemes for unknown distributions in a variety of scenarios; see, for example, Barron [17] and Barron and Cover [18]. It is also similar to notions of complexity in Dudley's entropy integral (cf. Dudley [77]) in empirical process theory, where the fluctuations of an empirical process are governed by a tradeoff between covering number and approximation of individual terms in the process.

10.2.2 Minimax bounds using global packings

There is now a four step process to proving minimax lower bounds using the global Fano method. Our starting point is to recall the Fano minimax lower bound in Proposition 9.4.3, which begins with the construction of a set of points $\{\theta(P_v)\}_{v \in \mathcal{V}}$ that form a 2δ -packing of a set Θ in some ρ -semimetric. With this inequality in mind, we perform the following four steps:

- (i) *Bound the packing entropy.* Give a lower bound on the packing number of the set Θ with 2δ -separation (call this lower bound $M(\delta)$).
- (ii) *Bound the metric entropy.* Give an upper bound on the KL-metric entropy of the class \mathcal{P} of distributions containing all the distributions P_v , that is, an upper bound on $\log N_{\text{kl}}(\epsilon, \mathcal{P})$.
- (iii) *Find the critical radius.* Noting as in Corollary 10.2.2 that with n i.i.d. observations, we have

$$I(V; X_1, \dots, X_n) \leq \inf_{\epsilon \geq 0} \{n\epsilon^2 + \log N_{\text{kl}}(\epsilon, \mathcal{P})\},$$

we now balance the information $I(V; X_1^n)$ and the packing entropy $\log M(\delta)$. To that end, we choose ϵ_n and $\delta > 0$ at the *critical radius*, defined as follows: choose the any ϵ_n such that

$$n\epsilon_n^2 \geq \log N_{\text{kl}}(\epsilon_n, \mathcal{P}),$$

and choose the largest $\delta_n > 0$ such that

$$\log M(\delta_n) \geq 4n\epsilon_n^2 + 2\log 2 \geq 2N_{\text{kl}}(\epsilon_n, \mathcal{P}) + 2n\epsilon_n^2 + 2\log 2 \geq 2(I(V; X_1^n) + \log 2).$$

(We could have chosen the ϵ_n attaining the infimum in the mutual information, but this way we need only an upper bound on $\log N_{\text{kl}}(\epsilon, \mathcal{P})$.)

- (iv) *Apply the Fano minimax bound.* Having chosen δ_n and ϵ_n as above, we immediately obtain that for the Markov chain $V \rightarrow X_1^n \rightarrow \hat{V}$,

$$\mathbb{P}(V \neq \hat{V}) \geq 1 - \frac{I(V; X_1, \dots, X_n) + \log 2}{\log M(\delta_n)} \geq 1 - \frac{1}{2} = \frac{1}{2},$$

and thus, applying the Fano minimax bound in Proposition 9.4.3, we obtain

$$\mathfrak{M}_n(\theta(\mathcal{P}); \Phi \circ \rho) \geq \frac{1}{2}\Phi(\delta_n).$$

10.2.3 Example: non-parametric regression

In this section, we flesh out the outline in the prequel to show how to obtain a minimax lower bound for a non-parametric regression problem directly with packing and metric entropies. In this example, we sketch the result, leaving explicit constant calculations to the dedicated reader. Nonetheless, we recover an analogue of Theorem 10.1.4 on minimax risks for estimation of 1-Lipschitz functions on $[0, 1]$.

We use the standard non-parametric regression setting, where our observations Y_i follow the independent noise model (10.1.1), that is, $Y_i = f(X_i) + \varepsilon_i$. Letting

$$\mathcal{F} := \{f : [0, 1] \rightarrow \mathbb{R}, f(0) = 0, f \text{ is Lipschitz}\}$$

be the family of 1-Lipschitz functions with $f(0) = 0$, we have

Proposition 10.2.3. *There exists a universal constant $c > 0$ such that*

$$\mathfrak{M}_n(\mathcal{F}, \|\cdot\|_\infty) := \inf_{\hat{f}_n} \sup_{f \in \mathcal{F}} \mathbb{E}_f \left[\|\hat{f}_n - f\|_\infty \right] \geq c \left(\frac{\sigma^2}{n} \right)^{1/3},$$

where \hat{f}_n is constructed based on the n independent observations $f(X_i) + \varepsilon_i$.

The rate in Proposition 10.2.3 is sharp to within factors logarithmic in n ; a more precise analysis of the upper and lower bounds on the minimax rate yields

$$\mathfrak{M}_n(\mathcal{F}, \|\cdot\|_\infty) := \inf_{\hat{f}_n} \sup_{f \in \mathcal{F}} \mathbb{E}_f \left[\|\hat{f}_n - f\|_\infty \right] \asymp \left(\frac{\sigma^2 \log n}{n} \right)^{1/3}.$$

See, for example, Tsybakov [182] for a proof of this fact.

Proof Our first step is to note that the covering and packing numbers of the set \mathcal{F} in the ℓ_∞ metric satisfy

$$\log N(\delta, \mathcal{F}, \|\cdot\|_\infty) \asymp \log M(\delta, \mathcal{F}, \|\cdot\|_\infty) \asymp \frac{1}{\delta}. \quad (10.2.4)$$

To see this, fix some $\delta \in (0, 1)$ and assume for simplicity that $1/\delta$ is an integer. Define the sets $E_j = [\delta(j-1), \delta j)$, and for each $v \in \{-1, 1\}^{1/\delta}$ define $h_v(x) = \sum_{j=1}^{1/\delta} v_j \mathbf{1}\{x \in E_j\}$. Then define the function $f_v(t) = \int_0^t h_v(t) dt$, which increases or decreases linearly on each interval of width δ in $[0, 1]$. Then these f_v form a 2δ -packing and a 2δ -cover of \mathcal{F} , and there are $2^{1/\delta}$ such f_v . Thus the asymptotic approximation (10.2.4) holds.

JCD Comment: TODO: Draw a picture

Now, if for some fixed $x \in [0, 1]$ and $f, g \in \mathcal{F}$ we define P_f and P_g to be the distributions of the observations $f(x) + \varepsilon$ or $g(x) + \varepsilon$, we have that

$$D_{\text{kl}}(P_f \| P_g) = \frac{1}{2\sigma^2} (f(X_i) - g(X_i))^2 \leq \frac{\|f - g\|_\infty^2}{2\sigma^2},$$

and if P_f^n is the distribution of the n observations $f(X_i) + \varepsilon_i$, $i = 1, \dots, n$, we also have

$$D_{\text{kl}}(P_f^n \| P_g^n) = \sum_{i=1}^n \frac{1}{2\sigma^2} (f(X_i) - g(X_i))^2 \leq \frac{n}{2\sigma^2} \|f - g\|_\infty^2.$$

In particular, this implies the upper bound

$$\log N_{\text{kl}}(\epsilon, \mathcal{P}) \lesssim \frac{1}{\sigma\epsilon}$$

on the KL-metric entropy of the class $\mathcal{P} = \{P_f : f \in \mathcal{F}\}$, as $\log N(\delta, \mathcal{F}, \|\cdot\|_\infty) \asymp \delta^{-1}$. Thus we have completed steps (i) and (ii) in our program above.

It remains to choose the critical radius in step (iii), but this is now relatively straightforward: by choosing $\epsilon_n \asymp (1/\sigma n)^{1/3}$, and whence $n\epsilon_n^2 \asymp (n/\sigma^2)^{1/3}$, we find that taking $\delta \asymp (\sigma^2/n)^{1/3}$ is sufficient to ensure that $\log N(\delta, \mathcal{F}, \|\cdot\|_\infty) \gtrsim \delta^{-1} \geq 4n\epsilon_n^2 + 2\log 2$. Thus we have

$$\mathfrak{M}_n(\mathcal{F}, \|\cdot\|_\infty) \gtrsim \delta_n \cdot \frac{1}{2} \gtrsim \left(\frac{\sigma^2}{n}\right)^{1/3}$$

as desired. □

JCD Comment: Should we do higher-dimensional stuff?

10.3 Strong converses and high-probability lower bounds

The results we have developed so far provide what we might call *weak* converse results: if one attempts estimate a parameter, there is a constant probability of error for estimating to within a certain accuracy. The information theory literature, on the other hand, provides *strong converses*, which for our purposes we interpret as follows: if one attempts to send too much information

through a communication channel, then the probability of error in decoding messages necessarily approaches 1. To set the stage, recall the setting of Section 9.4, where we choose $V \in \mathcal{V}$ uniformly at random, then have the Markov chain $V \rightarrow Y \rightarrow \hat{V}$. Then letting $\epsilon = \mathbb{P}(\hat{V} \neq V)$, Fano's inequality (Corollary 9.4.2) equivalently states that

$$\log \text{card}(\mathcal{V}) \leq \frac{I(V; Y)}{1 - \epsilon} + \frac{\log 2}{1 - \epsilon}.$$

In a communication setting, where we wish to send a message $v \in \mathcal{V}$ over a noisy channel, this result states that if we wish to have vanishing error $\epsilon \rightarrow 0$, then the maximum number of messages it is possible to send has $\log \text{card}(\mathcal{V}) \leq I(V; Y) + \log 2$.

An elegant way to derive strong converses is to provide refined versions of Fano's inequality. To develop these refinements, we typically consider somewhat more specific settings than the completely arbitrary Markov chain $V \rightarrow Y \rightarrow \hat{V}$. The most common scenarios are independent sampling scenarios, where conditional on $V = v$, we draw $Y_1^n \in \mathcal{Y}^n$ from a product distribution P_v^n . Letting ϵ be some measure of the probability of error and $g(\epsilon)$ be a function for which $g(\epsilon) < \infty$ whenever $\epsilon < 1$, then a typical refinement of Fano's inequality takes the form

$$\log \text{card}(\mathcal{V}) \leq I(V; Y_1^n) + \tilde{O}(\sqrt{n} \cdot g(\epsilon)), \quad (10.3.1)$$

where the big- \tilde{O} notation may hide logarithmic factors.

That inequality (10.3.1) is stronger than the standard Fano inequality may not be immediately apparent, so it is useful to consider an abstract communication scenario, where we are allowed to use a fixed communication channel n times. Here, we think of V as a message from a (large) collection \mathcal{V} of messages, and upon choosing a message v to send, we encode it into a vector $x \in \mathcal{X}^n$; the i th output of the channel Y_i is then drawn independently according to $P(\cdot | x_i)$. Graphically, Figure 10.3 shows the setting. The most common scenario is communication of a binary signal

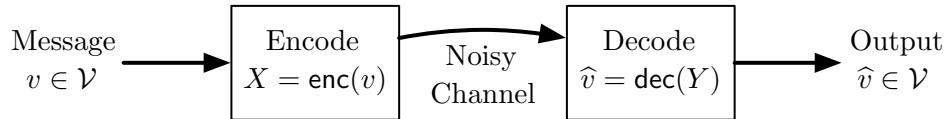


Figure 10.3: Communication over a noisy channel

where $\mathcal{X} = \mathcal{Y} = \{0, 1\}$, so that we send bits, and we think of the system as encoding a message $v \in \mathcal{V}$ as a bit string $x \in \{0, 1\}^n$, and the channel corrupts each bit $x_i \in \{0, 1\}$ independently to $Y_i \in \{0, 1\}$. Then we expect to be able to encode roughly 2^{cn} messages, for some constant c , and the information $I(V; Y_1^n)$ should grow linearly in n , as we send n messages. So inequality (10.3.1) would then imply we *must* have $\epsilon \rightarrow 1$ whenever $\log \text{card}(\mathcal{V}) - I(V; Y_1^n) \gg \sqrt{n}$, ignoring logarithmic factors.

We provide one exemplar result of the form (10.3.1) for finite output spaces \mathcal{Y} in the next section. We then show how it implies certain high-probability lower bounds for different estimation problems. In particular, we will show that for many estimation problems, there is some accuracy threshold ϵ_n for which the probability of estimating to accuracy *at all* better than ϵ_n tends to zero exponentially quickly.

10.3.1 Refined Fano inequalities

We begin by revisiting Fano's inequality, Corollary 9.4.2. We provide a proof of the inequality distinct from the “book” proof in Proposition 2.3.3, instead using the Donsker-Varadhan variational characterization of the KL-divergence, Theorem 6.1.1. The advantage of this approach is that the general idea immediately extends to allow us to use the blowing-up lemma (Corollary 7.3.2) to develop the refined Fano inequality.

For the first statement, we continue to work in the abstract scenario that we have a Markov chain $V \rightarrow Y \rightarrow \hat{V}$ without relying on any structure in the channel. For simplicity, we assume that the “decoder” $\hat{V} = \hat{v}(Y)$ is deterministic.

Proposition 10.3.1. *Let V be uniform on a set of size M and $\epsilon = \mathbb{P}(\hat{v}(Y) \neq V)$ be the average error of the estimator \hat{v} . Then*

$$(1 - \epsilon) \log(M + 1) \leq I(V; Y) + \log 2.$$

Proof When V is uniform, conditional on $V = v$, we draw $Y \sim P_v$. Let $Q = \frac{1}{M} \sum_v P_v$ be the marginal distribution on Y . Then by the definition of the mutual information and the Donsker-Varadhan variational inequality (Theorem 6.1.1), for any functions g_v we have

$$\begin{aligned} I(V; Y) &= \frac{1}{M} \sum_v D_{\text{kl}}(P_v \| Q) \geq \frac{1}{M} \sum_v \mathbb{E}_v[g_v] - \frac{1}{M} \sum_v \log \mathbb{E}_Q[e^{g_v}] \\ &\geq \frac{1}{M} \sum_v \mathbb{E}_v[g_v] - \log \left(\frac{1}{M} \sum_v \mathbb{E}_Q[e^{g_v}] \right), \end{aligned} \quad (10.3.2)$$

where the final inequality uses the concavity of the logarithm and \mathbb{E}_v denotes expectation with respect to P_v . Now we make a judicious choice of the function g_v . Let $A_v = \{y \in \mathcal{Y} \mid \hat{v}(y) = v\}$, so that the sets A_v partition \mathcal{Y} , and for a $\delta > 0$ to be chosen, define

$$g_v(y) = \log(\mathbf{1}\{y \in A_v\} + \delta).$$

Use the shorthands $p_v = P_v(A_v)$ and $q_v = Q(A_v)$, so that $\sum_v q_v = 1$. Then $\mathbb{E}_v[g_v] = p_v \log(1 + \delta) + (1 - p_v) \log \delta$, while $\mathbb{E}_Q[e^{g_v}] = \delta + q_v$, and so

$$I(V; Y) \geq \frac{1}{M} \sum_v \left[p_v \log(1 + \delta) - (1 - p_v) \log \frac{1}{\delta} \right] + \log M - \log(1 + M\delta).$$

Finally, we use the definition of the average error $\epsilon = \frac{1}{M} \sum_v (1 - p_v)$, and so

$$I(V; Y) \geq (1 - \epsilon) \log(1 + \delta) - \epsilon \log \frac{1}{\delta} + \log M - \log(1 + \delta M).$$

Choose $\delta = \frac{1}{M}$ and simplify. □

Proposition 10.3.1 recovers (with the ever-so-slightly stronger quantity $\log(M + 1)$ instead of $\log M$) the initial (weak) Fano inequality in Corollary 9.4.2. The coming theorem extends this result when the channel involves repeated independent sampling, where conditional on $V = v$, we draw a vector $Y_1^n \sim P_v^n$, where P_v^n is a product distribution on the output space \mathcal{Y}^n . As in the strong converse for hypothesis testing (Proposition 7.3.5), we provide bounds on the probabilities of “enlarged” subsets of a partition \mathcal{Y}^n . Note that in this variant, we control the *maximal* error rather than the average error present in the weaker Fano inequalities; some weakening like this is unavoidable (though the proof of that is beyond our scope).

Theorem 10.3.2. *Let V be uniform on a set of size M and $\epsilon = \max_v \mathbb{P}(\hat{v}(Y_1^n) = v \mid V = v)$ be the maximal error of the estimator \hat{v} . Then for all $t \geq \sqrt{8/n}$,*

$$(1 - e^{-t^2}) \log(M + 1) \leq I(V; Y_1^n) + \sqrt{\frac{n}{2}} \log(n \text{card}(\mathcal{Y})) \left[t + \sqrt{\log \frac{1}{1 - \epsilon}} \right]$$

Proof We mimic the proof of Proposition 10.3.1, beginning from the inequality (10.3.2). Here, however, instead of taking the functions g_v to be logarithms of the indicators of the partitions $A_v = \{y \in \mathcal{Y}^n \mid \hat{v}(y) = v\}$, we instead consider the r -blowups $A_{vr} := \{y \in \mathcal{Y}^n \mid d_{\text{ham}}(y, A_v) \leq r\}$, and define

$$g_v(y) := \log(\mathbf{1}\{y \in A_{vr}\} + \delta)$$

for some $\delta > 0$ to be chosen. Let

$$p_v^r := P_v^n(A_{vr}) \quad \text{and} \quad q_v^r := Q(A_{vr})$$

for shorthand. Then because $\mathbb{E}_v[g_v] = p_v^r \log(1 + \delta) - (1 - p_v^r) \log \frac{1}{\delta}$ and $\mathbb{E}_Q[e^{g_v}] = q_v^r + \delta$, inequality (10.3.2) implies

$$I(V; Y_1^n) \geq \frac{1}{M} \sum_v \left[p_v^r \log(1 + \delta) - (1 - p_v^r) \log \frac{1}{\delta} \right] + \log M - \log \left(\sum_v q_v^r + M\delta \right). \quad (10.3.3)$$

To leverage inequality (10.3.3), we must upper bound $\sum_v q_v^r$, which we rewrite as

$$\sum_v q_v^r = \sum_v Q(A_{vr}) = \sum_{y \in \mathcal{Y}^n} \sum_{v: y \in A_{vr}} Q(\{y\}).$$

For any fixed $y_0 \in \mathcal{Y}^n$, a quick counting argument yields

$$\text{card}\{y \in \mathcal{Y}^n \mid d_{\text{ham}}(y, y_0) \leq r\} \leq \text{card}(\mathcal{Y})^r \binom{n}{r}.$$

So because the sets $\{A_v\}$ partition \mathcal{Y}^n , any vector y belongs to at most $\binom{n}{r} \text{card}(\mathcal{Y})^r$ of the r -enlargements A_{vr} , and we have

$$\sum_v q_v^r \leq \sum_{y \in \mathcal{Y}^n} \binom{n}{r} \text{card}(\mathcal{Y})^r Q(\{y\}) = \binom{n}{r} \text{card}(\mathcal{Y})^r.$$

Substituting this into inequality (10.3.3) and letting $\epsilon(r) := \frac{1}{M} \sum_v (1 - p_v^r)$ be the blown-up average “error” we obtain

$$I(V; Y_1^n) \geq (1 - \epsilon(r)) \log(1 + \delta) - \epsilon(r) \log \frac{1}{\delta} + \log M - \log \left(\binom{n}{r} \text{card}(\mathcal{Y})^r + \delta M \right).$$

Now we use the blowing-up lemma (Corollary 7.3.2). For $t \geq 0$ define

$$r = r(t) := t \sqrt{\frac{n}{2}} + \sqrt{\frac{n}{2} \log \frac{1}{1 - \epsilon}},$$

so that for this r we have

$$p_v^r \geq 1 - e^{-t^2} \quad \text{and} \quad \epsilon(r) \leq e^{-t^2}.$$

Substituting this above and taking $\delta = 1/M$, we have

$$\begin{aligned} I(V; Y_1^n) &\geq \left(1 - e^{-t^2}\right) \log \left(1 + \frac{1}{M}\right) + \left(1 - e^{-t^2}\right) \log M - \log \left(1 + \binom{n}{r} \text{card}(\mathcal{Y})^r\right) \\ &= \left(1 - e^{-t^2}\right) \log(M+1) - \log \left(1 + \binom{n}{r} \text{card}(\mathcal{Y})^r\right). \end{aligned} \quad (10.3.4)$$

Recognizing that $1 + \binom{n}{r} \text{card}(\mathcal{Y})^r \leq (n \text{card}(\mathcal{Y}))^r$ for $r \geq 2$, inequality (10.3.4) implies

$$I(V; Y_1^n) \geq \left(1 - e^{-t^2}\right) \log(1+M) - r [\log n + \log \text{card}(\mathcal{Y})].$$

Substitute for $r = r(t)$ above and rearrange, recognizing it is sufficient that $t \geq \sqrt{8/n}$ to guarantee that $r \geq 2$. \square

The cardinality restriction on \mathcal{Y} in Theorem 10.3.2 is, while inelegant, not typically too onerous—for example, if we represent Y as a 32-bit floating point number then $\log \text{card}(\mathcal{Y}) \leq 32 \log 2$. For large enough n , we obtain a cleaner statement. Suppose that $\log n + \log \text{card}(\mathcal{Y}) \leq \sqrt{2} \log n$, and using $\frac{1}{\sqrt{2}-1} = 1 + \sqrt{2}$, it is sufficient that $n \geq \text{card}(\mathcal{Y})^{1+\sqrt{2}}$. Then for all $t \geq \sqrt{8/n}$,

$$\left(1 - e^{-t^2}\right) \log(M+1) \leq I(V; Y_1^n) + t \sqrt{n \log^2 n} + \sqrt{n \log^2 n \cdot \log \frac{1}{1-\epsilon}}. \quad (10.3.5)$$

We can also rearrange the refined Fano inequality in Theorem 10.3.2 to provide direct lower bounds on probabilities of error in communication and estimation settings.

Proposition 10.3.3. *Let the conditions of Theorem 10.3.2 hold and $n \geq \text{card}(\mathcal{Y})^{\sqrt{2}+1}$. Define*

$$\gamma_n := \frac{\log M - I(V; Y_1^n)}{\sqrt{n} \log n} - \sqrt{2 \log \log M} - \frac{1}{\sqrt{n} \log n}.$$

Then the maximal probability of error ϵ satisfies

$$\epsilon := \max_v P_v(\hat{V} \neq v) \geq 1 - \exp(-[\gamma_n]_+^2).$$

Proof We make the choice $t = \sqrt{2 \log \log M}$ in the bound (10.3.5). Then $e^{-t^2} \log(M+1) \leq \frac{\log(M+1)}{2 \log M} \leq 1$, and

$$\sqrt{n \log^2 n \log \frac{1}{1-\epsilon}} \geq \log M - I(V; Y_1^n) - \sqrt{2 \log \log M} \sqrt{n \log^2 n} - 1.$$

Divide both sides by $\sqrt{n} \log n$ to obtain

$$\sqrt{\log \frac{1}{1-\epsilon}} \geq \frac{\log M - I(V; Y_1^n)}{\sqrt{n} \log n} - \sqrt{2 \log \log M} - \frac{1}{\sqrt{n} \log n}.$$

Solve for ϵ . \square

Theorem 10.3.2 has the weakness that it relies explicitly on the finiteness of the output space \mathcal{Y} . Extensions of the result exist, though their proofs rely on reverse hypercontractivity of Markov semigroups and related functional inequalities and so are beyond our scope. We assume the same setting as Theorem 10.3.2, where $V \in \mathcal{V}$ is uniform, and conditional on $V = v$, we draw $Y_i \stackrel{\text{ind}}{\sim} P_{v,i}$ for $i = 1, \dots, n$, that is, $Y_1^n \sim P_v^n$ for a product distribution P_v^n . Then instead of a cardinality bound on \mathcal{Y} , we assume there exists a baseline probability measure P_0 on \mathcal{Y} providing the uniform likelihood ratio bound

$$\alpha := \max_i \max_v \left\| \frac{dP_{v,i}}{dP_0} \right\|_\infty < \infty.$$

When \mathcal{Y} is finite, we can always take P_0 to be uniform on \mathcal{Y} , so that $\alpha \leq \text{card}(\mathcal{Y})$. Then Liu et al. [139, Theorem 3.2] prove the following result:

Theorem 10.3.4. *Let $\hat{v} : \mathcal{Y}^n \rightarrow \mathcal{V}$ be an estimator of $V \sim \text{Uniform}(\mathcal{V})$, where $M = \text{card}(\mathcal{V})$. Let $\epsilon < 1$ and assume the geometric average of correct estimates probability satisfies*

$$\left(\prod_{v \in \mathcal{V}} \mathbb{P}(\hat{v}(Y_1^n) = v \mid V = v) \right)^{1/M} \geq 1 - \epsilon.$$

Then

$$\log M \leq I(V; Y_1^n) + 2\sqrt{n(\alpha - 1)} \sqrt{\log \frac{1}{1 - \epsilon}} + \log \frac{1}{1 - \epsilon}.$$

Theorem 10.3.4 exhibits a better dependence on n and the probability of error ϵ —using the geometric average rather than the maximum probability of error—than Theorem 10.3.2 and so asymptotically is stronger.

10.3.2 High probability estimation lower bounds

Connecting the refinement of Fano’s inequality in Theorem 10.3.2 to estimation problems provides new insights into the fundamental limits of estimation. As in our development of the local Fano method, we will work with parametric problems for which $D_{\text{kl}}(P_\theta \| P_{\theta'}) \lesssim \|\theta - \theta'\|_2^2$, but we will require slightly more delicate control. We consider probabilistic models for predicting a target y from a covariate vector $x \in \mathcal{X} \subset \mathbb{R}^d$, so that

$$Y \mid X = x \sim P_\theta(\cdot \mid x).$$

We say that the model has κ *quadratic information bound* if for each $x \in \mathcal{X}$, there exists a distribution $P_0(\cdot \mid x)$ on Y such that

$$D_{\text{kl}}(P_\theta(\cdot \mid x) \| P_0(\cdot \mid x)) \leq \kappa^2 (x^\top \theta)^2, \quad (10.3.6)$$

refining the local quadratic bound (9.4.6). Many models satisfy such bounds; typical generalized linear models satisfy it, for example (recall Chapter 3.4). Concretely, Example 3.4.8 shows that for binary logistic regression of a label $y \in \{0, 1\}$, the null model P_0 that $Y \sim \text{Uniform}\{0, 1\}$ satisfies

$$D_{\text{kl}}(P_\theta(\cdot \mid x) \| P_0(\cdot \mid x)) \leq \min \left\{ \frac{1}{8} (x^\top \theta)^2, \log 2 \right\}.$$

The key is that, for any model with a quadratic information bound (10.3.6), we can upper bound the information in estimation problems. We proceed conditionally on the covariates $X_1^n \in \mathcal{X}^n$, and take our packing set $\mathcal{V} \subset \{-1, 1\}^d$.

JCD Comment: Ideally, connect this to the reduction from estimation to testing and everything. Also, fix the constants a little bit.

Theorem 10.3.5. Fix a covariate matrix $X = [x_1 \cdots x_n]^\top \in \mathbb{R}^{n \times d}$, and assume the model P_θ satisfies the quadratic information bound (10.3.6) for $x \in \{x_i\}_{i=1}^n$. Then there exists a numerical constant $c > 0$ such that the following holds. For $\delta \geq 0$, define

$$\gamma(\delta) := \frac{cd - \kappa^2 \delta^2 \|X\|_{\text{Fr}}^2}{\sqrt{n} \log n} - \sqrt{2 \log 2 \cdot \log d} - \frac{1}{\sqrt{n} \log n}.$$

Then for all $n \geq \text{card}(\mathcal{Y})^{\sqrt{2}+1}$, there exists θ with $\|\theta\|_2 \leq \delta \sqrt{d}$ such that

$$\mathbb{P}_\theta \left(\|\hat{\theta}(Y_1^n) - \theta\|_2 \geq \frac{1}{2} \delta \sqrt{d} \right) \geq 1 - \exp(-[\gamma(\delta)]_+^2).$$

Unpacking Theorem 10.3.5, let us assume the covariates are standardized to $x_i \in \{\pm 1\}^d$, so that $\|X\|_{\text{Fr}}^2 = nd$. Then the constant $\gamma \geq \frac{cd - \kappa^2 \delta^2 nd}{\sqrt{n} \log n} - \sqrt{2 \log d}$, and taking $\delta^2 = \frac{c}{2n\kappa^2}$, we obtain

$$\gamma(\delta) \geq \frac{cd}{2\sqrt{n} \log n} - \sqrt{2 \log d}.$$

We thus obtain the following corollary, which shows that there is essentially no probability that an estimator can have accuracy better than $O(1)\sqrt{d/n}$ when the dimension scales so that $d^2/n \gg 1$.

Corollary 10.3.6. In addition to the conditions of Theorem 10.3.5, assume the covariates $x_i \in [-1, 1]^d$ and the dimension $d \geq \sqrt{n} \log^3 n$. Then there exists a numerical constant $c > 0$ such that

$$\sup_{\|\theta\|_2 \leq \kappa \sqrt{d/n}} \mathbb{P}_\theta \left(\|\hat{\theta}_n - \theta\|_2 \geq \frac{c}{\kappa} \sqrt{\frac{d}{n}} \right) \geq 1 - n^{-c}.$$

Of course, if the dimension is larger relative to n , we obtain stronger bounds; for example, once $d \geq \sqrt{n} \log^4 n$ we obtain

$$\sup_{\theta} \mathbb{P}_\theta \left(\|\hat{\theta}_n - \theta\|_2 \geq \frac{c}{\kappa} \sqrt{\frac{d}{n}} \right) \geq 1 - \frac{1}{n^{c \log n}}.$$

JCD Comment: Specialize to logistic regression for giggles. Commentary that can do better with bounded likelihood ratios. Maybe add an exercise to that effect? Also clean up the conditional on x part.
Clean up / connect conditions on $d^2/n \rightarrow \infty$ or whatever.

10.3.3 Proof of Theorem 10.3.5

As in our discussion of the local Fano method, we will construct a packing set $\mathcal{V} \subset \{-1, 1\}^d$, and for each $v \in \mathcal{V}$ identify $\theta_v = \delta v$ for some $\delta > 0$ to be chosen.

The key for us will be to construct a packing that has *both* strong separation and for which $V \sim \text{Uniform}(\mathcal{V})$ has small second moment matrix, in particular, satisfying $\mathbb{E}[VV^\top] \preceq (1 + o(1))I_d$. The probabilistic method allows us to do this:

Lemma 10.3.7. *There exist numerical constants $0 < c$ and $C < \infty$ such that a packing $\mathcal{V} \subset \{-1, 1\}^d$ of the hypercube exists with the following properties: its cardinality satisfies $\text{card}(\mathcal{V}) \geq \exp(cd)$,*

$$\|v - v'\|_2 > \sqrt{d} \text{ for } v \neq v' \in \mathcal{V}, \text{ and } \mathbb{E}[VV^\top] \preceq \left(1 + C\sqrt{d/e^{cd}}\right) I_d$$

for $V \sim \text{Uniform}(\mathcal{V})$.

We defer the proof of the lemma temporarily, as it is a straightforward application of the concentration guarantees we have already developed.

Now, consider the refined Fano bound in Theorem 10.3.2. Let \mathcal{V} be the packing Lemma 10.3.7 guarantees, and for a $\delta \geq 0$ to be chosen, set $\theta_v = \delta v$ for each $v \in \mathcal{V}$. Let $V \sim \text{Uniform}(\mathcal{V})$, and conditional on $V = v$, draw

$$Y_i \stackrel{\text{iid}}{\sim} P_{\theta_v}(\cdot \mid X = x_i).$$

Let $\hat{\theta}(Y_1^n)$ be any estimator of θ . Because θ_v is well-separated, with $\|\theta_v - \theta_{v'}\|_2 > \delta\sqrt{d}$ we see that if $\|\hat{\theta} - \theta_v\|_2 \leq \frac{\delta}{2}\sqrt{d}$, then $\|\hat{\theta} - \theta_{v'}\|_2 > \frac{\delta}{2}\sqrt{d}$ for all $v \neq v'$. Thus, the test

$$\Psi(Y_1^n) := \begin{cases} v & \text{if } \|\hat{\theta} - \theta_v\|_2 \leq \frac{\delta}{2}\sqrt{d} \\ \text{arbitrary} & \text{otherwise} \end{cases}$$

satisfies

$$\mathbb{P}(\Psi(Y_1^n) \neq v \mid V = v) \leq \mathbb{P}\left(\|\hat{\theta} - \theta_v\|_2 > \frac{\delta}{2}\sqrt{d} \mid V = v\right).$$

Proposition 10.3.3 then implies that the maximal probability of error $\epsilon := \max_v \mathbb{P}(\|\hat{\theta} - \theta_v\|_2 > \frac{\delta}{2}\sqrt{d})$ satisfies $\epsilon \geq 1 - \exp(-[\gamma]_+^2)$ for

$$\gamma := \frac{\log M - I(V; Y_1^n)}{\sqrt{n} \log n} - \sqrt{2 \log \log M} - \frac{1}{\sqrt{n} \log n}$$

where $2^d \geq M \geq \exp(cd)$.

For the final step, we upper bound the mutual information. Letting $\bar{P}^n = \frac{1}{\text{card}(\mathcal{V})} \sum_{v \in \mathcal{V}} P_v^n$ and P_0^n be any other product distribution, we obtain

$$I(V; Y_1^n) = \frac{1}{\text{card}(\mathcal{V})} \sum_{v \in \mathcal{V}} D_{\text{kl}}(P_v^n \parallel \bar{P}^n) \leq \frac{1}{\text{card}(\mathcal{V})} \sum_{v \in \mathcal{V}} D_{\text{kl}}(P_v^n \parallel P_0^n).$$

Taking P_0 to be the assumed conditional distribution providing the quadratic bound (10.3.6), we have $D_{\text{kl}}(P_v^n \parallel P_0^n) \leq \kappa^2 \sum_{i=1}^n (x_i^\top \theta_v)^2$, and so using $\theta_v = \delta v$,

$$I(V; Y_1^n) \leq \kappa^2 \delta^2 \mathbb{E} \left[\sum_{i=1}^n (x_i^\top V)^2 \right] \leq \kappa^2 \delta^2 \left(1 + C\sqrt{d/e^{cd}}\right) \sum_{i=1}^n \|x_i\|_2^2$$

by Lemma 10.3.7. Substituting in γ above gives the theorem.

Proof of Lemma 10.3.7 Let the vectors $U_i \stackrel{\text{iid}}{\sim} \text{Uniform}\{\pm 1\}^d$ be uniform on the binary hypercube, $i = 1, \dots, N$, where we will choose N presently. Then for $i \neq j$ $\|U_i - U_j\|_2^2 = 2d - 2\langle U_i, U_j \rangle$, and $\langle U_i, U_j \rangle \stackrel{\text{dist}}{=} \langle \mathbf{1}_d, U_i \rangle$ by independence. So for $t \geq 0$,

$$\mathbb{P}(\|U_i - U_j\|_2^2 \leq 2d - 2t) = \mathbb{P}(\langle \mathbf{1}, U_i \rangle \geq t) \leq \exp\left(-\frac{t^2}{2d}\right)$$

by standard sub-Gaussian concentration. By a union bound, we have

$$\mathbb{P}\left(\min_{i \neq j} \|U_i - U_j\|_2^2 \leq 2d - 2t\right) \leq \binom{N}{2} \exp\left(-\frac{t^2}{2d}\right) \leq \exp\left(-\frac{t^2}{2d} + 2 \log N\right),$$

and the choice $t = d/2$ gives $\mathbb{P}(\min_{i \neq j} \|U_i - U_j\|_2^2 \leq d) \leq \exp(-d/8 + 2 \log N)$. Taking $N = \exp(\frac{d}{16} - \frac{1}{2})$ gives that $\|U_i - U_j\|_2 > \sqrt{d}$ for all $i \neq j$ with probability at least $1 - 1/e$.

To obtain a covariance bound, we use either Proposition 5.1.11 or Proposition 7.1.5, which for $N \geq \exp(cd)$ guarantees that $\frac{1}{N} \sum_{i=1}^N U_i U_i^\top \preceq (1 + C\sqrt{(d+t)/N})I_d$ with probability at least $1 - e^{-t}$. \square

JCD Comment: Here, put a strong converse for channel coding. Also put a strong converse for estimation in a $d^2/n \rightarrow \infty$ model for, e.g., logistic regression. Is it possible that the paper by Liu, van Handel, and Verdu gives better than the $d^2/n \log n$ convergence?

10.4 Exercises

JCD Comment: Some things to either do as exercises or include in the actual note

1. Density estimation (upper bound in 1-d)
2. Curse of dimensionality in nonparametric regression
3. Better version of the estimation lower bounds using Theorem 10.3.4.

Chapter 11

Constrained risk inequalities

In this chapter, we revisit our minimax bounds in the context of what we term *constrained risk inequalities*. While the minimax risk provides a first approach for providing fundamental limits on procedures, its reliance on the collection of *all* measurable functions as its class of potential estimators is somewhat limiting. Indeed, in most statistical and statistical learning problems, we have some type of constraint on our procedures: they must be efficiently computable, they must work with data arriving in a sequential stream, they must be robust, or they must protect the privacy of the providers of the data. In modern computational hardware, where physical limits prevent increasing clock speeds, we may like to use as much parallel computation as possible, though there are potential tradeoffs between “sequentialness” of procedures and their parallelism.

With this as context, we replace the minimax risk of Chapter 9.1 with the *constrained minimax risk*, which, given a collection \mathcal{C} of possible procedures—private, communication limited, or otherwise—defines

$$\mathfrak{M}(\theta(\mathcal{P}), \Phi \circ \rho, \mathcal{C}) := \inf_{\hat{\theta} \in \mathcal{C}} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\Phi(\rho(\hat{\theta}(X), \theta(P))) \right], \quad (11.0.1)$$

where as in the original defining equation (9.1.1) of the minimax risk, $\Phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is a nondecreasing loss, ρ is a semimetric on the space Θ , and the expectation is taken over the sample $X \sim P$. In this chapter, we study the quantity (11.0.1) via a few examples, highlighting possibilities and challenges with its analysis. We will focus on a restricted class of examples—many procedures do not fall in the framework we consider—that assumes, given a sample X_1, \dots, X_n , we can represent the class \mathcal{C} of estimators under consideration as acting on some view or processed version Z_i of X_i . This allows us to study communication complexity, memory complexity, and certain private estimators.

11.1 Strong data processing inequalities

The starting point for our results is to consider *strong data processing inequalities*, which improve upon the standard data processing inequality for divergences, as in Chapter 2.1.3, to provide more quantitative versions. The initial setting is straightforward: we have distributions P_0 and P_1 on a space \mathcal{X} , and a channel (Markov kernel) Q from \mathcal{X} to \mathcal{Z} . When Q is contractive on the space of distributions, we have a strong data processing inequality.

Definition 11.1 (Strong data processing inequalities). *Let $f : \mathbb{R}_+ \rightarrow \mathbb{R} \cup \{+\infty\}$ be convex and satisfy $f(1) = 0$. For distributions P_0, P_1 on \mathcal{X} and a channel Q from \mathcal{X} to a space \mathcal{Z} , define*

the marginal distribution $M_v(A) := \int Q(A \mid x) dP_v(x)$. The channel Q satisfies a strong data processing inequality with constant $\alpha \leq 1$ for the given f -divergence

$$D_f(M_0 \| M_1) \leq \alpha D_f(P_0 \| P_1)$$

for any choice of P_0, P_1 on \mathcal{X} . For any such f , we define the f -strong data processing constant

$$\alpha_f(Q) := \sup_{P_0 \neq P_1} \frac{D_f(M_0 \| M_1)}{D_f(P_0 \| P_1)}.$$

These types of inequalities are common throughout information and probability theory. Perhaps their most frequent use is in the development conditions for the fast mixing of Markov chains. Indeed, suppose the Markov kernel Q satisfies a strong data processing inequality with constant α with respect to variation distance. If π denotes the stationary distribution of the Markov kernel Q and we use the operator \circ to denote one step of the Markov kernel,¹

$$Q \circ P := \int Q(\cdot \mid x) dP(x),$$

then for any initial distribution π_0 on the space \mathcal{X} we have

$$\| \underbrace{Q \circ \dots \circ Q}_{k \text{ times}} \pi_0 - \pi \|_{\text{TV}} \leq \alpha^k \|\pi_0 - \pi\|_{\text{TV}}$$

because $Q \circ \pi = \pi$ by definition of the stationary distribution. Thus, the Markov chain enjoys geometric mixing.

To that end, a common quantity of interest is the *Dobrushin* coefficient, which immediately implies mixing rates.

Definition 11.2. The Dobrushin coefficient of a channel or Markov kernel Q is

$$\alpha_{\text{TV}}(Q) := \sup_{x, y} \|Q(\cdot \mid x) - Q(\cdot \mid y)\|_{\text{TV}}.$$

The Dobrushin coefficient satisfies many properties, some of which we discuss in the exercises and others of which we enumerate here. The first is that

Proposition 11.1.1. The Dobrushin coefficient is the strong data processing constant for the variation distance, that is,

$$\alpha_{\text{TV}}(Q) = \sup_{P_0 \neq P_1} \frac{\|Q \circ P_0 - Q \circ P_1\|_{\text{TV}}}{\|P_0 - P_1\|_{\text{TV}}}.$$

Proof There are two directions to the proof; one easy and one more challenging. For the easy direction, we see immediately that if $\mathbf{1}_x$ and $\mathbf{1}_y$ denote point masses at x and y , then

$$\sup_{P_0 \neq P_1} \frac{\|Q \circ P_0 - Q \circ P_1\|_{\text{TV}}}{\|P_0 - P_1\|_{\text{TV}}} \geq \sup_{x, y} \|Q(\cdot \mid x) - Q(\cdot \mid y)\|_{\text{TV}}$$

as $\|\mathbf{1}_x - \mathbf{1}_y\|_{\text{TV}} = 1$ for $x \neq y$.

¹The standard notation is usually to right-multiply the measure P , so that the marginal distribution $M = PQ$ means $M(A) = \int Q(A \mid x) dP(x)$; we find our notation more intuitive.

The other direction—that $\|Q \circ P_0 - Q \circ P_1\|_{\text{TV}} \leq \alpha_{\text{TV}} \|P_0 - P_1\|_{\text{TV}}$ —is more challenging. For this, recall Lemma 2.2.4 characterizing the variation distance, and let $Q_\star(A) := \inf_y Q(A | y)$. Then by definition of the Dobrushin coefficient $\alpha = \alpha_{\text{TV}}(Q)$, we evidently have $|Q(A | x) - Q_\star(A)| \leq \alpha$. Let $M_v = \int Q(\cdot | x) dP_v(x)$ for $v \in \{0, 1\}$. By expanding $dP_0 - dP_1$ into its positive and negative parts, we thus obtain

$$\begin{aligned} M_0(A) - M_1(A) &= \int Q(A | x) (dP_0 - dP_1)(x) \\ &= \int Q(A | x) [dP_0(x) - dP_1(x)]_+ - \int Q(A | x) [dP_1(x) - dP_0(x)]_+ \\ &\leq \int Q(A | x) [dP_0(x) - dP_1(x)]_+ - \int Q_\star(A) [dP_1(x) - dP_0(x)]_+ \\ &= \int Q(A | x) [dP_0(x) - dP_1(x)]_+ - \int Q_\star(A) [dP_0(x) - dP_1(x)]_+, \end{aligned}$$

where the final equality uses Lemma 2.2.4. But of course we then obtain

$$M_0(A) - M_1(A) = \int (Q(A | x) - Q_\star(A)) [dP_0(x) - dP_1(x)]_+ \leq \alpha \int [dP_0 - dP_1]_+ = \alpha \|P_0 - P_1\|_{\text{TV}},$$

where the inequality follows as $0 \leq Q(A | x) - Q_\star(A) \leq \alpha$ and the equality is one of the characterizations of the total variation distance in Lemma 2.2.4. \square

A more substantial fact is that the Dobrushin coefficient upper bounds *every* other strong data processing constant.

Theorem 11.1.2. *Let $f : \mathbb{R}_+ \rightarrow \mathbb{R} \cup \{\infty\}$ satisfy $f(1) = 0$. Then for any channel Q ,*

$$\alpha_{\text{TV}}(Q) \geq \alpha_f(Q).$$

The theorem is roughly a consequence of a few facts. First, Proposition 11.1.1 holds. Second, without loss of generality we may assume that $f \geq 0$; indeed, replace $f(t)$ with $h(t) = f(t) - f'(1)t$ for any $f'(1) \in \partial f(1)$, we have $h \geq 0$ as $0 \in \partial h(1)$ and $D_h = D_f$. Third, any $f \geq 0$ with $0 \in \partial f(1)$ can be approximated arbitrarily accurately with functions of the form $h(t) = \sum_{i=1}^k a_i [t - c_i]_+ + \sum_{i=1}^k b_i [d_i - t]_+$, where $c_i \geq 1$ and $d_i \leq 1$. For such h , an argument shows that

$$D_h(Q \circ P_0 \| Q \circ P_1) \leq \alpha_{\text{TV}}(Q) D_h(P_0 \| P_1),$$

which follows from the similarities between variation distance, with $f(t) = \frac{1}{2}|t|$, and the positive part functions $[\cdot]_+$.

There is a related result, which we do not prove, that guarantees that strong data processing constants for χ^2 -divergences are the “worst” constants. In particular, if $QP = \int Q(\cdot | x) dP(x)$ denotes the application of one step of a channel Q to $X \sim P$, then the χ^2 contraction coefficient is

$$\alpha_{\chi^2}(Q) = \sup_{P_0 \neq P_1} \frac{D_{\chi^2}(QP_0 \| QP_1)}{D_{\chi^2}(P_0 \| P_1)}.$$

Then it is possible to show that for any twice continuously differentiable f on \mathbb{R}_{++} with $f''(1) > 0$,

$$\alpha_{\chi^2}(Q) \leq \alpha_f(Q), \tag{11.1.1}$$

and we also have $\alpha_{\chi^2}(Q) = \alpha_{\text{kl}}(Q)$, so that the strong data processing inequalities for KL-divergence and χ^2 -divergence coincide.

In our context, that of (constrained) minimax lower bounds, such data processing inequalities immediately imply somewhat sharper lower bounds than the (unconstrained) applications in previous chapters. Indeed, let us revisit the situation present in the local Fano bound, where we the KL divergence has a Euclidean structure as in the bound (9.4.6), meaning that $D_{\text{kl}}(P_0 \| P_1) \leq \kappa^2 \delta^2$ when our parameters of interest $\theta_v = \theta(P_v)$ satisfy $\rho(\theta_0, \theta_1) \leq \delta$. We assume that the constraints \mathcal{C} impose that the data X_i is passed through a channel Q with KL-data processing constant $\alpha_{\text{KL}}(Q) \leq 1$. In this case, in the basic Le Cam's method (9.3.2), an application of Pinsker's inequality yields that whenever $\rho(\theta_0, \theta_1) \geq 2\delta$ then

$$\mathfrak{M}_n(\theta(\mathcal{P}), \Phi \circ \rho, \mathcal{C}) \geq \frac{\Phi(\delta)}{2} \left[1 - \sqrt{\frac{n}{2} D_{\text{kl}}(M_0 \| M_1)} \right] \geq \frac{\Phi(\delta)}{2} \left[1 - \sqrt{n \kappa^2 \alpha_{\text{KL}}(Q) \delta^2 / 2} \right],$$

and the “standard” choice of δ to make the probability of error constant results in $\delta^2 = (2n\kappa^2\alpha_{\text{KL}}(Q))^{-1}$, or the minimax lower bound

$$\mathfrak{M}_n(\theta(\mathcal{P}), \Phi \circ \rho, \mathcal{C}) \geq \frac{1}{4} \Phi \left(\frac{1}{\sqrt{2n\kappa^2\alpha_{\text{KL}}(Q)}} \right),$$

which suggests an effective sample size degradation of $n \mapsto n\alpha_{\text{KL}}(Q)$. Similarly, in the local Fano method in Chapter 9.4.1, we see identical behavior and an effective sample size degradation of $n \mapsto n\alpha_{\text{KL}}(Q)$, that is, if without constraints a sample size of $n(\epsilon)$ is required to achieve some desired accuracy ϵ , with the constraint a sample size of at least $n(\epsilon)/\alpha_{\text{KL}}(Q)$ is necessary.

11.2 Local privacy

In Chapter 8 on differential privacy, we define *locally private mechanisms* (Definition 8.2) as those for which there is no trust: individuals randomize their own data, and no central curator collects or analyzes and then privatizes the resulting statistics. With such privacy mechanisms, we can directly develop strong data processing inequalities, after which we can prove strong lower bounds on estimation. In this section, we (more or less) focus on one-dimensional quantities and Le Cam's two-point method for lower bounds, as they allow the most direct application of the ideas. We will later develop more sophisticated techniques.

We begin with our setting. We have a ϵ -differentially private channel Q taking inputs $x \in \mathcal{X}$ and outputting Z . Here, we allow *sequential interactivity*, meaning that the i th private variable Z_i may depend on both X_i and Z_1^{i-1} (see the graphical model in Figure 11.1), so that instead of the basic constraint in Definition 8.2 that $Q(A | x) \leq e^\epsilon Q(A | x')$ for all x, x' , local differential privacy instead means

$$\frac{Q(Z_i \in A | X_i = x, z_1^{i-1})}{Q(Z_i \in A | X_i = x', z_1^{i-1})} \leq e^\epsilon \quad (11.2.1)$$

for all (measurable) sets A and inputs x, x', z_1^{i-1} . The key result is the following contraction inequality on the space of probabilities.

Theorem 11.2.1. *Let Q be an ϵ -locally differentially private channel from \mathcal{X} to \mathcal{Z} . Then for any distributions P_0, P_1 inducing marginal distributions $M_v(\cdot) = \int Q(\cdot | x) dP_v(x)$,*

$$D_{\text{kl}}(M_0 \| M_1) + D_{\text{kl}}(M_1 \| M_0) \leq 4(e^\epsilon - 1)^2 \|P_0 - P_1\|_{\text{TV}}^2.$$

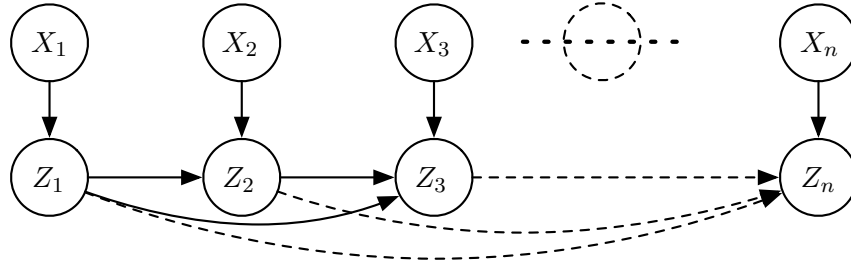


Figure 11.1. The sequentially interactive private observation model: the i th output Z_i may depend on X_i and the previously released Z_1^{i-1} .

Proof Without loss of generality, we assume that the output space \mathcal{Z} is finite (by definition (2.2.3)), and let $m_v(z)$ and $q(z | x)$ be the p.m.f.s of M and Q , respectively, and let P_0 and P_1 have densities p_0 and p_1 with respect to a measure μ . Then

$$D_{\text{kl}}(M_0 \| M_1) + D_{\text{kl}}(M_1 \| M_0) = \sum_z (m_0(z) - m_1(z)) \log \frac{m_0(z)}{m_1(z)}$$

For any $a, b \geq 0$, we have $\log \frac{a}{b} = \log(1 + \frac{a}{b} - 1) \leq \frac{a}{b} - 1$, and similarly, $\log \frac{b}{a} \leq \frac{b}{a} - 1$. That is, $|\log \frac{a}{b}| \leq \frac{|a-b|}{\min\{a,b\}}$. Substituting above, we obtain

$$D_{\text{kl}}(M_0 \| M_1) + D_{\text{kl}}(M_1 \| M_0) \leq \sum_z \frac{(m_0(z) - m_1(z))^2}{\min\{m_0(z), m_1(z)\}}.$$

To control the difference $m_0(z) - m_1(z)$, note that for any fixed $x_0 \in \mathcal{X}$ we have

$$\int_{\mathcal{X}} q(z | x_0) (p_0(x) - p_1(x)) d\mu(x) = 0.$$

Thus

$$m_0(z) - m_1(z) = \int_{\mathcal{X}} (q(z | x) - q(z | x_0)) (p_0(x) - p_1(x)) d\mu(x),$$

and so

$$\begin{aligned} |m_0(z) - m_1(z)| &\leq \sup_{x \in \mathcal{X}} |q(z | x) - q(z | x_0)| \int_{\mathcal{X}} |p_0(x) - p_1(x)| d\mu(x) \\ &= 2q(z | x_0) \sup_{x \in \mathcal{X}} \left(\frac{q(z | x)}{q(z | x_0)} - 1 \right) \|P_0 - P_1\|_{\text{TV}}. \end{aligned}$$

By definition of local differential privacy, $\frac{q(z|x)}{q(z|x_0)} - 1 \leq e^\epsilon - 1$, and as x_0 was arbitrary we obtain

$$|m_0(z) - m_1(z)| \leq 2(e^\epsilon - 1) \inf_x q(z | x) \|P_0 - P_1\|_{\text{TV}}.$$

Noting that $\inf_x q(z | x) \leq \min\{m_0(z), m_1(z)\}$ we obtain the theorem. \square

To be able to apply this result to obtain minimax lower bounds for estimation as in Section 9.3, we need to address samples drawn from product distributions, even with the potential interaction (11.2.1). In this case, we consider sequential samples $Z_i \sim Q(\cdot \mid X_i, Z_1^{i-1})$ and define $M_v^n = \int Q(\cdot \mid x_1^n) dP_v(x_1^n)$ to be the marginal distribution over all the Z_1^n . Then we have the following corollary.

Corollary 11.2.2. *Assume that each channel $Q(\cdot \mid X_i, Z_1^{i-1})$ is ε_i -differentially private. Then*

$$D_{\text{kl}}(M_0^n \parallel M_1^n) \leq 4 \sum_{i=1}^n (e^{\varepsilon_i} - 1)^2 \|P_0 - P_1\|_{\text{TV}}^2.$$

Proof Recalling the chain rule (2.1.6) for the KL-divergence, we have

$$D_{\text{kl}}(M_0^n \parallel M_1^n) = \sum_{i=1}^n \mathbb{E}_{M_0} [D_{\text{kl}}(M_{0,i}(\cdot \mid Z_1^{i-1}) \parallel M_{1,i}(\cdot \mid Z_1^{i-1}))],$$

where the outer expectation is taken over Z_1^{i-1} drawn marginally from M_0^n , and $M_{v,i}(\cdot \mid z_1^{i-1})$ denotes the conditional distribution on Z_i given $Z_1^{i-1} = z_1^{i-1}$ when $X_1^n \stackrel{\text{iid}}{\sim} P_v$. Writing this distribution out, we note that Z_i is conditionally independent of $X_{\setminus i}$ given X_i and Z_1^{i-1} by construction, so for any set A

$$\begin{aligned} M_{v,i}(A \mid z_1^{i-1}) &= \int Q(Z_i \in A \mid x_1^n, z_1^{i-1}) dP_v(x_1^n \mid z_1^{i-1}) = \int Q(Z_i \in A \mid x_i, z_1^{i-1}) dP_v(x_1^n \mid z_1^{i-1}) \\ &= \int Q(Z_i \in A \mid x_i, z_1^{i-1}) dP_v(x_i). \end{aligned}$$

Now we know that $Q(Z_i \in \cdot \mid x_i, z_1^{i-1})$ is ε_i -differentially private by assumption, so Theorem 11.2.1 gives

$$D_{\text{kl}}(M_{0,i}(\cdot \mid z_1^{i-1}) \parallel M_{1,i}(\cdot \mid z_1^{i-1})) \leq 4(e^{\varepsilon_i} - 1)^2 \|P_0 - P_1\|_{\text{TV}}^2$$

for any realization z_1^{i-1} of Z_1^{i-1} . Iterating this gives the result. \square

Local privacy is such a strong condition on the channel Q that it actually “transforms” the KL-divergence into a variation distance, so that even if two distributions P_0 and P_1 have infinite KL-divergence $D_{\text{kl}}(P_0 \parallel P_1) = +\infty$ —for example, if their supports are not completely overlapping—their induced marginals have the much smaller divergence $D_{\text{kl}}(M_0 \parallel M_1) \leq 4(e^\varepsilon - 1)^2 \|P_0 - P_1\|_{\text{TV}}^2 \leq 4(e^\varepsilon - 1)^2$. This transformation into a different metric means that even in estimation problems that should on their faces be easy become quite challenging under local privacy constraints; for example, minimax squared error for estimating the mean of a random variable with finite variance scales as $1/\sqrt{n}$ rather than the typical $1/n$ scaling in non-private cases (see Exercise 11.4).

Let us demonstrate how to apply Corollary 11.2.2 in a few applications. Our main object of interest is the private analogue of the minimax risk (9.1.1), where for a parameter $\theta : \mathcal{P} \rightarrow \Theta$, semimetric ρ , and loss Φ , for a family of channels \mathcal{Q} we define the *channel-constrained minimax risk*

$$\mathfrak{M}_n(\theta(\mathcal{P}), \Phi \circ \rho, \mathcal{Q}) := \inf_{\hat{\theta}_n} \inf_{Q \in \mathcal{Q}} \sup_{P \in \mathcal{P}} \mathbb{E}_{P,Q} [\Phi(\rho(\hat{\theta}_n(Z_1^n), \theta(P)))] . \quad (11.2.2)$$

When we take $\mathcal{Q} = \mathcal{Q}_\varepsilon$ to be the collection of ε -locally differentially private (interactive) channels (11.2.1), we obtain the ε -locally private minimax risk.

A few examples showing lower (and upper) bounds for the private minimax risk (11.2.2) in mean estimation follow.

Example 11.2.3 (Bounded mean estimation): Let \mathcal{P} be the collection of distributions with supports on $[-b, b]$, where $0 < b < \infty$. Then for any $\varepsilon \geq 0$, the minimax squared error satisfies

$$\mathfrak{M}_n(\theta(\mathcal{P}), (\cdot)^2, \mathcal{Q}_\varepsilon) \gtrsim \frac{b^2}{(e^\varepsilon - 1)^2 n} + \frac{b^2}{n}.$$

The second term in the bound is the classic minimax rate for this collection of distributions. To see the first term, take Bernoulli distributions P_0 and $P_1 \in \mathcal{P}$, where for some $\delta \geq 0$ to be chosen, under P_0 we have $X = b$ with probability $\frac{1-\delta}{2}$ and $-b$ otherwise, while under P_1 we have $X = b$ with probability $\frac{1+\delta}{2}$ and $X = -b$ otherwise. Then $\|P_0 - P_1\|_{\text{TV}} = \delta$, $\mathbb{E}_1[X] - \mathbb{E}_0[X] = 2b\delta$, and by Le Cam's method (9.3.3), for any ε -locally private channel Q and induced marginals M_0^n, M_1^n as in Corollary 11.2.2, we have

$$\begin{aligned} \mathfrak{M}_n(\theta(\mathcal{P}), (\cdot)^2, \{Q\}) &\geq \frac{b^2 \delta^2}{2} \left(1 - \sqrt{\frac{1}{2} D_{\text{kl}}(M_0^n \| M_1^n)} \right) \geq \frac{b^2 \delta^2}{2} \left(1 - \sqrt{2(e^\varepsilon - 1)^2 n \|P_0 - P_1\|_{\text{TV}}^2} \right) \\ &= \frac{b^2 \delta^2}{2} \left(1 - \sqrt{2(e^\varepsilon - 1)^2 n \delta^2} \right). \end{aligned}$$

Setting $\delta^2 = \frac{1}{8n(e^\varepsilon - 1)^2}$ gives the claimed minimax bound. \diamond

Effectively, then, we see a reduction in the effective sample size: when ε is large, there is no change, but otherwise, the estimation error is similar to that when we observe a sample of size $n\varepsilon^2$.

Example 11.2.4 (Estimating the parameter of a uniform distribution): In exercise 9.2, we show that estimating the parameter θ of a $\text{Uniform}(\theta, \theta + 1)$ distribution has minimax squared error scaling as $1/n^2$. Under local differential privacy, this is impossible. Let $\mathcal{P} = \{\text{Uniform}(\theta, \theta + 1), \theta \in [0, 1]\}$ be the collection of uniform distributions with the given supports. Letting P_0 and P_1 be $\text{Uniform}(0, 1)$ and $\text{Uniform}(\delta, 1 + \delta)$, respectively, where $\delta \geq 0$ is to be chosen, we have $\|P_0 - P_1\|_{\text{TV}} = \delta$, while for any ε -differentially private channel Q and induced marginals M_0 and M_1 ,

$$D_{\text{kl}}(M_0^n \| M_1^n) \leq 4(e^\varepsilon - 1)^2 n \|P_0 - P_1\|_{\text{TV}}^2 = 4(e^\varepsilon - 1)^2 n \delta^2.$$

Applying Le Cam's method (9.3.3) and taking $\delta \asymp \frac{1}{\sqrt{n(e^\varepsilon - 1)}}$, we thus have that if \mathcal{Q}_ε denotes the collection of ε -locally differentially private channels,

$$\mathfrak{M}_n(\theta(\mathcal{P}), (\cdot)^2, \mathcal{Q}_\varepsilon) \gtrsim \frac{1}{(e^\varepsilon - 1)^2 n}.$$

When $\varepsilon \lesssim 1$, the best attainable rate thus scales as $\frac{1}{n\varepsilon^2}$. \diamond

In both the preceding examples, a number of simple estimators achieve the given minimax rates. The simplest is one based on the Laplace mechanism (Example 8.1.3): let $W_i \stackrel{\text{iid}}{\sim} \text{Laplace}(1)$, and set $Z_i = X_i + \frac{2b}{\varepsilon} W_i$ in Example 11.2.3 and $Z_i = X_i + \frac{2}{\varepsilon} W_i$ in Example 11.2.4. In the former, define $\hat{\theta}_n = \bar{Z}_n$ to be the mean; in the latter, $\mathbb{E}[\bar{Z}_n] = \frac{\theta+1}{2}$, so $\hat{\theta}_n = 2\bar{Z}_n - 1$ achieves the minimax rate.

More extreme examples are possible. Consider, for example, the problem of testing the support of a distribution, where we care only about distinguishing two distributions.

Example 11.2.5 (Support testing): Consider the problem of testing between the support of two uniform distributions, that is, given n observations, we wish to test whether $P = P_0 = \text{Uniform}[0, 1]$ or $P = P_1 = \text{Uniform}[\theta, 1]$ for some $\theta \in (0, 1)$. We can ask the rate at which we may take $\theta \downarrow 0$ with n while still achieving non-trivial testing power. Without privacy, a simple (and optimal) test Ψ is to simply check whether any observation $X_i < \theta$, in which case we can trivially accept P_0 and reject P_1 , otherwise accepting P_1 . Then

$$P_0(X_i > \theta, \text{ all } i) = (1 - \theta)^n \quad \text{while} \quad P_1(X_i > \theta, \text{ all } i) = 1.$$

So the summed probability of error

$$P_0(\Psi = 1) + P_1(\Psi = 0) = (1 - \theta)^n \leq \exp(-\theta n),$$

and if $\theta \gg 1/n$ this tends to zero, while $\theta_n = \theta_0/n$ yields $\lim_n P_0(\Psi = 1) = e^{-\theta_0}$.

Consider now the private case. Then for any ε -differentially private channel Q and induced marginals M_0, M_1 , we have $D_{\text{kl}}(M_0^n \| M_1^n) \leq 4n(e^\varepsilon - 1)^2 \|P_0 - P_1\|_{\text{TV}}^2$ by Corollary 11.2.2 while $\|P_0 - P_1\|_{\text{TV}} = \theta$. The Bretagnolle-Huber inequality (Proposition 2.2.8.(b)) thus guarantees that

$$\|M_0^n - M_1^n\|_{\text{TV}}^2 \leq 1 - \exp(-D_{\text{kl}}(M_0^n \| M_1^n)) \leq 1 - \exp(-4n(e^\varepsilon - 1)^2 \theta^2).$$

Whenever $\theta \ll \frac{1}{\sqrt{n}}$, we have $\|M_0^n - M_1^n\|_{\text{TV}} \rightarrow 0$, and so for *any* test based on the private data Z_1^n , the probabilities of error

$$\inf_{\Psi} \{P_0(\Psi(Z_1^n) \neq 0) + P_1(\Psi(Z_1^n) \neq 1)\} \geq 1 - \sqrt{1 - \exp(-c_\varepsilon n \theta^2)},$$

where $c_\varepsilon = 4(e^\varepsilon - 1)^2$. In the range that $\frac{1}{n} \ll \theta \ll \frac{1}{\sqrt{n}}$, then, there is an essentially exponential gap between the non-private and private cases. \diamond

11.3 Communication complexity

Communication complexity is a broad field, encompassing results establishing fundamental limits in streaming and online algorithms, memory-limited procedures, and (of course) in minimal communication in various fields. Recent connections between communication complexity and information-theoretic techniques have increased its applicability in statistical problems, which is our main motivation here, and to which we return in force in Section 11.4 to come. To motivate our approaches, however, we give a (necessarily limited) overview of communication complexity, along with some of the basic techniques and approaches, which then extend to statistical problems.

11.3.1 Classical communication complexity problems

The most basic problems in communication complexity are not really statistical, instead asking a simpler question: two entities (always named Alice and Bob) have inputs x, y and wish to jointly compute a function $f(x, y)$. The question is then how many bits—or other messages—Alice and Bob need to communicate to compute this value. Less abstractly, Alice and Bob have input domains \mathcal{X} and \mathcal{Y} (often, these are $\{0, 1\}^n$), and Alice receives a vector $x \in \mathcal{X}$ and Bob $y \in \mathcal{Y}$, each unknown to the other, and they jointly exchange messages until they can successfully evaluate $f(x, y)$. To abstract away any details of the computational model, we assume each has infinite computational

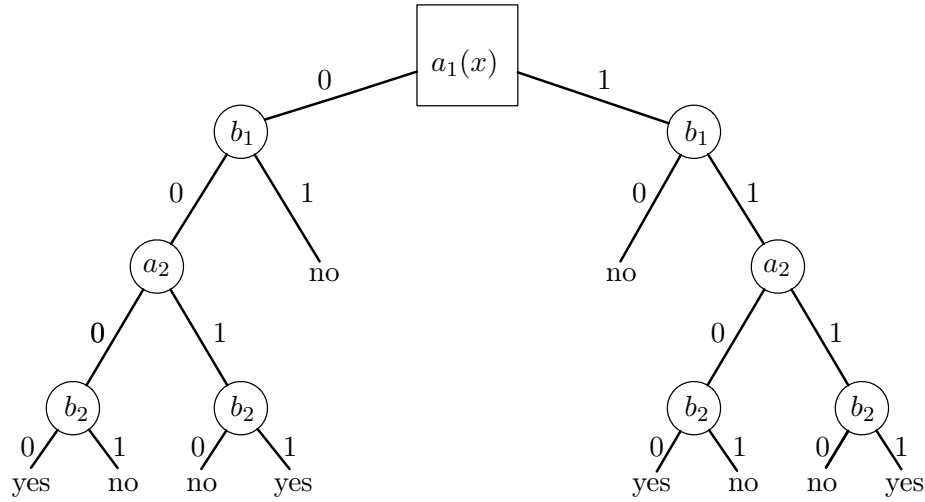


Figure 11.2. A communication tree representing testing equality for 2-dimensional bit strings $x, y \in \{0, 1\}^2$. Internal nodes labeled a_j communicate the j th bit $a_j(x) = x_j$ of x , while internal nodes labeled b_j communicate the j th bit $b_j(y) = y_j$ of y . The maximum number of messages is 4. (A more efficient protocol is to have Alice send the entire string $x \in \{0, 1\}^n$, then for Bob to check equality $x = y$ and output “Yes” or “No.”)

power, which allows a focus on communication. To formulate this as communication, we consider a *protocol* Π , which specifies the messages that each of Alice and Bob send to one another. We view this as a series of rounds, where at each round, the protocol allows one $\{0, 1\}$ -valued bit to be sent and determines who sends this bit, and, at termination time, can compute $f(x, y)$ based on the communicated message. Then the communication cost of Π is the maximum number of messages sent to (correctly) compute f over all inputs x, y .

A more convenient formulation for analysis is to consider a binary tree:

Definition 11.3. A protocol Π over a domain $\mathcal{X} \times \mathcal{Y}$ with output space \mathcal{Z} is a binary tree, where each internal node v is labeled with a mapping $a_v : \mathcal{X} \rightarrow \{0, 1\}$ or $b_v : \mathcal{Y} \rightarrow \{0, 1\}$ and each leaf is labeled with a value $z \in \mathcal{Z}$.

Then to execute a communication protocol Π on input (x, y) , we walk down the tree: beginning at the root node, for each internal node v labeled a_v (an Alice node) we walk left if $a_v(x) = 0$ and right if $a_v(x) = 1$, and each node v labeled b_v (a Bob node) we walk left if $b_v(y) = 0$ and right if $b_v(y) = 1$. Then the *communication cost* of the protocol Π is the height of the tree, which we denote by $\text{depth}(\Pi)$. Figure 11.2 shows an example for testing the equality $x = y$ of two 2-dimensional bit strings $x, y \in \{0, 1\}^2$.

In classical communication complexity, the main questions center around the *communication complexity* of a function $f : \mathcal{X} \rightarrow \mathcal{Y}$, which is the length of the shortest protocol that computes f correctly on all inputs: letting $\Pi_{\text{out}}(x, y)$ denote the final output of the protocol Π on inputs (x, y) , this is

$$\text{CC}(f) := \inf \{ \text{depth}(\Pi) \mid \Pi_{\text{out}}(x, y) = f(x, y) \text{ for all } x \in \mathcal{X}, y \in \mathcal{Y} \}.$$

In many cases, it is useful to allow randomized communication protocols, which tolerate some probability of error; in this case, we let Alice and Bob each have access to (an arbitrary amount) of randomness, which we can identify without loss of generality with uniform random variables

$U_a, U_b \stackrel{\text{iid}}{\sim} \text{Uniform}[0, 1]$, and the nodes a_v and b_v in Definition 11.3 are then mappings $a_v : \mathcal{X} \times [0, 1] \rightarrow \{0, 1\}$ and $b_v : \mathcal{Y} \times [0, 1] \rightarrow \{0, 1\}$ and they calculate $a_v(\cdot, U_a)$ and $b_v(\cdot, U_b)$, respectively. Abusing notation slightly by leaving this randomness implicit, the *randomized communication complexity* for an accuracy δ is then the length of the shortest randomized protocol that calculates $f(x, y)$ correctly with probability at least $1 - \delta$, that is,

$$\text{RCC}_\delta(f) := \inf \{ \text{depth}(\Pi) \mid \mathbb{P}(\Pi_{\text{out}}(x, y) \neq f(x, y)) \leq \delta \text{ for all } x \in \mathcal{X}, y \in \mathcal{Y} \}. \quad (11.3.1)$$

In the definition (11.3.1), we leave the randomization in Π implicit, and note that we require that the tree it induces still have a maximum length. We note that essentially any choice of $\delta > 0$ is immaterial, as we always have

$$\text{RCC}_\delta(f) \leq O(1) \log \frac{1}{\delta} \cdot \text{RCC}_{1/3}(f),$$

making all (constant) probability of error complexities essentially equivalent. (See Exercise 11.7.)

There are variants of randomized complexity that allow public randomness rather than private randomness, which can yield simpler algorithms and somewhat reduced complexity, but this improvement is limited, as Alice and Bob can always essentially simulate public randomness (see Exercise 11.8). Letting $\mathfrak{P}_{\text{pub}}$ be the collection of protocols in which both Alice and Bob have access to a shared random variable $U \sim \text{Uniform}[0, 1]$, we make the obvious extension

$$\text{RCC}_\delta^{\text{pub}}(f) := \inf_{\Pi \in \mathfrak{P}_{\text{pub}}} \{ \text{depth}(\Pi) \mid \mathbb{P}(\Pi_{\text{out}}(x, y, U) \neq f(x, y)) \leq \delta \text{ for all } x \in \mathcal{X}, y \in \mathcal{Y} \}.$$

Finally, we have *distributional* communication complexity, which for a probability measure μ on inputs $\mathcal{X} \times \mathcal{Y}$ is the depth of the shortest protocol that succeeds with a given μ -probability:

$$\text{DCC}_\delta^\mu(f) := \inf \{ \text{depth}(\Pi) \mid \mu(\Pi_{\text{out}}(X, Y) \neq f(X, Y)) \leq \delta \}, \quad (11.3.2)$$

where the infimum is taken over *deterministic* protocols.

The final notion we consider is the *information complexity*. In this case, we require again that for each input pair x, y , the (potentially randomized) protocol $\Pi(x, y)$ still compute $f(x, y)$ correctly with probability at least $1 - \delta$, but instead of measuring the depth of the tree, we let X, Y be drawn randomly from some distribution and measure the mutual information $I_2(X, Y; \Pi(X, Y))$. (We use base-2 logarithms to reflect bit communication.) In this case, we define

$$\text{IC}_\delta(f) := \sup_{\Pi} \inf \{ I_2(X, Y; \Pi(X, Y)) \mid \mathbb{P}(\Pi_{\text{out}}(x, y) \neq f(x, y)) \leq \delta \text{ for all } x \in \mathcal{X}, y \in \mathcal{Y} \}, \quad (11.3.3)$$

where the supremum is taken over joint distributions on (X, Y) , the infimum over randomized protocols Π , and the right probability \mathbb{P} is over any randomness in Π . There is a subtlety in this definition: we require Π to be accurate on *all* inputs (x, y) , not just with probability over the distribution on (X, Y) in the information measure $I(X, Y; \Pi(X, Y))$. Relaxations to distributional variants of the information complexity (11.3.3) are also natural, as in the definition (11.3.2). Thus we sometimes consider the distributional information complexity

$$\text{IC}_\delta^\mu(f) := \inf_{\Pi} \{ I_2(X, Y; \Pi(X, Y)) \mid \mu(\Pi_{\text{out}}(X, Y) \neq f(X, Y)) \leq \delta \},$$

where the infimum can be taken over deterministic or randomized protocols.

The different notions of communication complexity satisfy a natural ordering, making proving lower bounds for some notions (or conversely, developing low-communication methods for different protocols) much easier or harder than others. We record the standard inequalities in the coming proposition, which essentially follows immediately from the operational interpretation of entropy as the average length of the best encoding of a signal (Section 2.4.1).

Proposition 11.3.1. *For any function f , $\delta \in (0, 1)$, and probability measure μ on $\mathcal{X} \times \mathcal{Y}$,*

$$CC(f) \geq RCC_\delta(f) \geq RCC_\delta^{\text{pub}}(f) \geq DCC_\delta^\mu(f) \geq IC_\delta^\mu(f)$$

and

$$RCC_\delta(f) \geq IC_\delta(f).$$

Proof The first two inequalities are immediate. By Theorem 2.4.3, we have

$$\text{depth}(\Pi) \geq H_2(\Pi) \geq H_2(\Pi) - H_2(\Pi \mid X, Y) = I_2(X, Y; \Pi(X, Y)),$$

and so for all $\delta \in (0, \frac{1}{2})$ we have both

$$RCC_\delta(f) \geq IC_\delta(f) \quad \text{and} \quad DCC_\delta^\mu(f) \geq IC_\delta^\mu(f).$$

All that remains is to demonstrate $RCC_\delta^{\text{pub}}(f) \geq DCC_\delta^\mu(f)$. For this, let Π be any protocol with public randomness U such that $\mathbb{P}(\Pi_{\text{out}}(x, y, U) \neq f(x, y)) \leq \delta$ for all x, y . Then by taking an expectation over $(X, Y) \sim \mu$, we obtain

$$\delta \geq \mathbb{E}_\mu [\mathbb{P}(\Pi_{\text{out}}(X, Y, U) \neq f(X, Y) \mid X, Y)] \geq \inf_u \mu(\Pi_{\text{out}}(X, Y, u) \neq f(X, Y)),$$

that is, there must be at least some u achieving the average error of Π , and the protocol Π is deterministic given u . So any protocol Π using public randomness to achieve probability of error δ can be modified into a deterministic protocol $\Pi(\cdot, \cdot, u)$ that achieves μ -probability of error δ .² \square

Frequently, the first inequality in Proposition 11.3.1 is strict—even exponentially large—while the randomized complexity and information complexity end up being of roughly the same order. Understanding these differences is one of the major goals in communication complexity research.

11.3.2 Deterministic communication: lower bounds and structure

Deterministic communication complexity lower bounds often admit fairly elegant and somewhat elementary arguments, and the gaps between them and the randomized complexity highlight that we indeed expect providing lower bounds on randomized communication (11.3.1) or information (11.3.3) complexity to be quite challenging. The starting point, to which we will return when we consider randomized protocols, is to understand some structural aspects of the inputs and outputs of a protocol tree.

Recall that a set $R \subset \mathcal{X} \times \mathcal{Y}$ is a *rectangle* if it has the form $R = A \times B$ for some $A \subset \mathcal{X}$ and $B \subset \mathcal{Y}$. Equivalently, R is a rectangle if $(x_0, y_0) \in R$ and $(x_1, y_1) \in R$ imply that $(x_0, y_1) \in R$. As the next proposition shows, rectangular sets provide a key way to understand communication complexity.

²This is one direction of Yao's minimax theorem [193], which states that communication complexity with public (shared) randomness and worst-case distributional complexity are identical: $RCC_\delta^{\text{pub}}(f) = \sup_\mu DCC_\delta^\mu(f)$.

Proposition 11.3.2. *Let v be a node in a deterministic protocol Π and R_v be those pairs (x, y) reaching node v . Then R_v is a rectangle.*

Proof We prove the result by induction. Certainly, for the root node v , we have $R_v = \mathcal{X} \times \mathcal{Y}$, which is a rectangle. Now, let v be an arbitrary (non-root) node in the tree and w its parent; assume w.l.o.g. that v is the left child of w and that in w , Alice speaks (that is, we use $a_w : \mathcal{X} \rightarrow \{0, 1\}$.) Then $R_w = A \times B$ by the inductive assumption. If $a_w(x) = 0$, then

$$R_v = \{\{x\} \times B \mid a_w(x) = 0, x \in A\} = \{\{x \mid a_w(x) = 0\} \cap A\} \cap B,$$

which is a rectangle. \square

The structure of rectangles for correct protocols thus naturally determines the communication complexity of a function f . For a set $R \subset \mathcal{X} \times \mathcal{Y}$, we say R is f -constant if $f(x, y) = f(x', y')$ for all $(x, y) \in R$ and $(x', y') \in R$. Thus, any correct protocol Π necessarily partitions $\mathcal{X} \times \mathcal{Y}$ into a collection of f -constant rectangles, where we identify the rectangles with the leaves l of the protocol tree. In particular, Proposition 11.3.2 implies the following corollary.

Corollary 11.3.3. *Let N be the size of the minimal partition of $\mathcal{X} \times \mathcal{Y}$ into f -constant rectangles. Then $\text{CC}(f) \geq \log_2 N$.*

Proof Any correct protocol Π partitions $\mathcal{X} \times \mathcal{Y}$ into the f -constant rectangles $\{R_l\}$ indexed by its leaves l . The minimal depth of a binary tree with at least N leaves is $\log_2 N$. \square

A related corollary follows by considering *fooling sets*, which are basically sets that rectangles cannot contain.

Definition 11.4 (Fooling sets). *A set $S \subset \mathcal{X} \times \mathcal{Y}$ is a fooling set for f if for any two pairs $(x_0, y_0) \in S$ and $(x_1, y_1) \in S$ satisfying $f(x_0, y_0) = f(x_1, y_1)$, at least one of the inequalities $f(x_0, y_1) \neq f(x_0, y_0)$ or $f(x_1, y_0) \neq f(x_0, y_0)$ holds.*

With this definition, the next corollary is almost immediate.

Corollary 11.3.4. *Let f have a fooling set S of size N . Then $\text{CC}(f) \geq \log_2 N$.*

Proof By definition, no f -constant rectangle contains more than a single element of S . So the tree associated with any correct protocol Π has a single leaf for each element of S . \square

An extension of the fooling set idea is the *rectangle measure* method, which proves that (for some probability measure P) the “size” of f -constant rectangles is small. By judicious choice of the probability, we can then demonstrate lower bounds.

Proposition 11.3.5. *Let P be a probability distribution on $\mathcal{X} \times \mathcal{Y}$. If all f -constant rectangles R have probability at most $P(R) \leq \delta$, Then $\text{CC}(f) \geq \log_2 \frac{1}{\delta}$.*

Proof By the union bound, any f -constant partition of $\mathcal{X} \times \mathcal{Y}$ into rectangles $\{R_l\}_{l=1}^N$ satisfies $1 \leq \sum_{l=1}^N P(R_l) \leq N\delta$. So $N \geq \frac{1}{\delta}$, and the result follows by Corollary 11.3.3. \square

With these results, we can provide lower bounds on two exemplar problems that will inform much of our coming development.

Example 11.3.6 (Equality): Consider the problem of testing equality of two n -bit strings $x, y \in \{0, 1\}^n$, letting $f = \text{EQ}$ be $f(x, y) = 1$ if $x = y$ and 0 otherwise. Define the set $S = \{(x, x) \mid x \in \{0, 1\}^n\}$, which has cardinality 2^n , and satisfies $f(x, x) = 1$ for all $(x, x) \in S$. That S is a fooling set is immediate: for any (x, x) and $(x', x') \in S$, if $x \neq x'$, then certainly $(x, x') \notin S$. So

$$n \leq \text{CC}(\text{EQ}) \leq n + 1,$$

where the upper bound follows by letting Alice simply communicate the string x and Bob check if $x = y$, outputting 1 or 0 as $x = y$ or $x \neq y$. \diamond

The second example concerns inner products on \mathbb{F}_2 , the field of arithmetic on the integers modulo 2 (that is, with bit strings); one could extend this to inner products in more complicated number systems (such as floating point), but the basic ideas are cleaner when we deal with bits.

Example 11.3.7 (Inner products on \mathbb{F}_2): Consider computing the inner product $\text{IP}_2(x, y) = \langle x, y \rangle \bmod 2$ for n -bit strings $x, y \in \{0, 1\}^n$, where addition is performed modulo 2. Rather than constructing a fooling set directly, we use Proposition 11.3.5 and let P be the uniform distribution on $\{0, 1\}^n \times \{0, 1\}^n$. Let $R = A \times B$ be a rectangle with $\langle x, y \rangle = 0$ for all $x \in A$ and $y \in B$. The linearity of the inner product guarantees that $\langle x, y \rangle = 0$ for all $x \in \text{span}(A)$ and $y \in \text{span}(B)$, the (linear) spans of A and B in \mathbb{F}_2^n , respectively. Now recognize that $\text{span}(A), \text{span}(B) \subset \mathbb{F}_2^n$ are orthogonal subspaces of \mathbb{F}_2^n , and so their dimensions $d_0 = \dim(\text{span}(A))$ and $d_1 = \dim(\text{span}(B))$ satisfy $d_0 + d_1 \leq n$.

Noting that if $d_0 = \dim(A)$ then $|A| \leq 2^{d_0}$ in \mathbb{F}_2^n , we thus obtain $|R| \leq |A| \cdot |B| \leq 2^n$, which (under the uniform measure P) satisfies

$$P(R) \leq \frac{2^n}{2^{2n}} = 2^{-n}.$$

By Proposition 11.3.5, we thus have

$$n \leq \text{CC}(\text{IP}_2) \leq n + 1,$$

where once again the upper bound follows by letting Alice simply communicate $x \in \{0, 1\}^n$ and having Bob output $\langle x, y \rangle \bmod 2$. \diamond

11.3.3 Randomization, information complexity, and direct sums

When we allow randomization, the complexity bounds can, in some cases, drastically change. Consider again the equality function in Example 11.3.6. When we allow randomization, we can achieve $O(\log n)$ complexity to check equality (with high probability).

Example 11.3.8 (Equality with randomization): Let $x, y \in \{0, 1\}^n$ and p be a prime number satisfying $n^2 \leq p \leq 2n^2$ (the Prime Number Theorem guarantees the existence of such a p). Let Alice choose a uniformly random number $U \in \{0, \dots, p-1\}$ and compute the polynomial

$$a(U) = x_1 + x_2U + x_3U^2 + \dots + x_nU^{n-1} \bmod p.$$

Then Alice may communicate both U and $a(U)$ to Bob, which requires at most $2 \log_2 p \leq 4 \log_2 n + 2 \log 2$ bits. Then Bob checks whether

$$b(U) = y_1 + y_2U + y_3U^2 + \dots + y_nU^{n-1} \bmod p$$

satisfies $b(U) = a(U)$. If so, Bob outputs “Yes” (equality), and otherwise, Bob outputs “No.” This protocol satisfies $\text{depth}(\Pi) \leq 4 \log_2 n + 1$. Moreover, if $x = y$, it is always correct, while if $x \neq y$, then the protocol is incorrect only if $a(U) = b(U)$, that is, U is a root of the polynomial

$$p(u) = \sum_{i=1}^n (x_i - y_i) u^{i-1}.$$

But this is a non-zero degree $n - 1$ polynomial, which has at most $n - 1$ roots (on the field \mathbb{F}_p ; see Appendix A.1 for a brief review of polynomials). Thus for $x \neq y$ we have

$$\mathbb{P}(\Pi(x, y) \text{ fails}) = \mathbb{P}(a(U) = b(U)) \leq \frac{n-1}{p} < \frac{1}{n},$$

and so $\text{RCC}_{1/n}(\text{EQ}) \leq O(1) \log n$, exponentially improving over deterministic complexity.

In passing, we make two additional remarks. First, this protocol is one-way and non-interactive: Alice can simply send $O(\log n)$ bits. Second, we can achieve essentially any probability of success in the bound while still only paying logarithmically in communication, as taking $n^k \leq p \leq 2n^k$ for $k \geq 2$ yields $\text{RCC}_{1/n^k}(\text{EQ}) \leq 2k \log_2 n + O(1)$. \diamond

Example 11.3.8 makes clear that any lower bounds on randomized communication complexity, or, relatedly, information complexity, will necessarily be somewhat more subtle than those we have presented for CC. We develop a few of the main ideas here. Because our focus is on information theoretic techniques, we pass over a few of the standard tools for proving lower bounds involving *discrepancy* and randomized inputs, touching on these in the bibliographic notes at the end of the chapter. One of our main goals will be to show that the information complexity of the inner product is indeed $\Omega(n)$, a much stronger result than Example 11.3.7. In contrast to the lower bounds we provide for minimax risk in most of this book, the focus in communication complexity is to take an *a priori* accurate estimator and demonstrate that it *requires* a certain amount of information to be communicated, rather than the contrapositive result that limited information yields inaccurate estimators. While these are clearly equivalent, it can be fruitful to use the perspective most relevant for the problem at hand.

Two main ideas form the basis for information complexity lower bounds: first, *direct sum* inequalities, which show that computing a function on n inputs requires roughly order n more communication than computing it (or at least, one of the constituent functions making it up) on one. The second important insight is to provide lower bounds on the information necessary to compute different primitives, and the particular structure of even randomized communication protocols makes this possible. For the remainder of Section 11.3.3, we address the first of these, returning to the information complexity of primitives in Section 11.3.4.

Direct sum bounds and decomposition

To show direct sum inequalities, we demonstrate that computing some function on n inputs requires roughly n times the communication of single-input computation. In general, we consider functions f of the form

$$f(x_1^n, y_1^n) = g(h(x_1, y_1), h(x_2, y_2), \dots, h(x_n, y_n)), \quad (11.3.4)$$

where g is the global function of the n primitives h , calling such functions decomposable with primitive h . Several problems have the decomposable structure (11.3.4); focusing on the case that the inputs $x, y \in \{0, 1\}^n$ and $f(x, y) \in \{0, 1\}$, we have the following three immediate examples.

Example 11.3.9 (Composition in equality): The equality function $f(x, y) = 1$ if $x \neq y$ and $f(x, y) = 0$ otherwise satisfies the decomposition (11.3.4), where $h(x_i, y_i) = \mathbf{1}\{x_i \neq y_i\}$ and g is the OR function $g(z) = \mathbf{1}\{\langle \mathbf{1}, z \rangle > 0\}$, which is 1 if any of z_1, \dots, z_n is non-zero, and 0 otherwise. \diamond

Example 11.3.10 (Decomposition of inner product): The inner product in \mathbb{F}_2 , $f(x, y) = \langle x, y \rangle \bmod 2$, where $h(x_i, y_i) = x_i y_i$, and $g(z) = \langle \mathbf{1}, z \rangle \bmod 2$, which satisfies $g(z) = 0$ if $\sum_{i=1}^n z_i$ is even and $g(z) = 1$ otherwise. \diamond

Example 11.3.11 (Decomposition of disjointness): The set disjointness function $f(x, y) = \text{DISJ}(x, y) := \mathbf{1}\{\langle x, y \rangle > 0\}$ arises when x, y are characteristic vectors of two subsets A, B of $[n]$, that is, $x_i = \mathbf{1}\{i \in A\}$ and $y_i = \mathbf{1}\{i \in B\}$. Then $f(x, y) = \mathbf{1}\{A \cap B \neq \emptyset\}$, which corresponds to g being the OR $g(z) = \mathbf{1}\{\langle \mathbf{1}, z \rangle > 0\}$ and h the AND function $h(x_i, y_i) = x_i y_i$. \diamond

While Example 11.3.8 makes clear that the decomposition (11.3.4) is not sufficient to guarantee a randomized complexity lower bound of order n , it will be useful.

To develop the main information complexity direct sum theorem showing that the information complexity of f is at least the sum of the complexities of its constituent primitives, we leverage what we term *plantable inputs*:

Definition 11.5. Let $f : \mathcal{X}^n \times \mathcal{Y}^n \rightarrow \{0, 1\}$ have the decomposition (11.3.4), where the primitive h is $\{0, 1\}$ -valued. The pair $(x, y) \in \mathcal{X}^n \times \mathcal{Y}^n$ admits a planted solution if for each $i \in \{1, \dots, n\}$, all x'_i, y'_i , and vectors all

$$x' = (x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n) \quad \text{and} \quad y' = (y_1, \dots, y_{i-1}, y'_i, y_{i+1}, \dots, y_n),$$

we have $f(x', y') = h(x'_i, y'_i)$.

The binary inner product in Examples 11.3.7 and 11.3.10 has many plantable inputs: any of the 3^n pairs of vectors $x, y \in \{0, 1\}^n$ with $\langle x, y \rangle = 0$ admit planted solutions, as we have $x_i y_i = 0$ for each i . The set-disjointness problem, Example 11.3.11, has the same plantable inputs. For the equality function, only the 2^n pairs $x = y$ admit planted solutions.

We outline the key idea to our direct sum lower bounds. Because we define information complexity for protocols Π that are correct on all inputs with high probability, we can choose an arbitrary distribution on inputs $(x_1^n, y_1^n) \in \mathcal{X}^n \times \mathcal{Y}^n$. Thus we choose a *fooling distribution* μ for f , meaning that for $(X_i, Y_i) \stackrel{\text{iid}}{\sim} \mu$ the pair $(X_1^n, Y_1^n) \in \mathcal{X}^n \times \mathcal{Y}^n$ always admits a planted solution (Definition 11.5). The next definition says this slightly differently.

Definition 11.6. A distribution μ on $(x, y) \in \mathcal{X} \times \mathcal{Y}$ is a fooling distribution if all (x_1^n, y_1^n) in the support of the product μ^n admit planted solutions (Definition 11.5).

Typically, fooling distributions μ require some dependence between X_i and Y_i —for example, in the inner product, we require $X_i Y_i = 0$, so that if $X_i = 1$ then $Y_i = 0$ and vice versa:

Example 11.3.12 (A fooling distribution for inner products and set disjointness): Define the distribution μ on pairs $(x, y) \in \{0, 1\} \times \{0, 1\}$ as follows: let V be uniform on $\{0, 1\}$, and conditional on $V = 0$, set $X = 0$ and let $Y \sim \text{Uniform}\{0, 1\}$; conditional on $V = 1$, set $Y = 0$ and let $X \sim \text{Uniform}\{0, 1\}$. Then certainly $XY = 0$, and any set of pairs $(X_i, Y_i) \stackrel{\text{iid}}{\sim} \mu$ satisfy both that the binary inner product $\text{IP}_2(X_1^n, Y_1^n) = \langle X_1^n, Y_1^n \rangle \bmod 2 = 0$ and set disjointness $\text{DISJ}(X_1^n, Y_1^n) = \mathbf{1}\{\langle X_1^n, Y_1^n \rangle > 0\} = 0$. \diamond

Fooling distributions, as in Example 11.3.12, make conditioning natural in information complexity. If $(X, Y) \sim \mu$, there is always a random variable V such that $X \perp\!\!\!\perp Y \mid V$, that is, X and Y are conditionally independent given V (trivially, we can take $V = X$). Thus, for function $h : \mathcal{X} \times \mathcal{Y} \rightarrow \{0, 1\}$, we define the *conditional information complexity*

$$\text{CIC}_\delta^\mu(h) := \inf_{\Pi} \sup_V \{I(X, Y; \Pi(X, Y) \mid V) \text{ s.t. } \mathbb{P}(\Pi_{\text{out}}(x, y) \neq h(x, y)) \leq \delta \text{ for all } x \in \mathcal{X}, y \in \mathcal{Y}\},$$

where the infimum is over all (randomized) protocols and the supremum is over all random variables making X and Y conditionally independent with joint distribution $(X, Y) \sim \mu$. So if we can find a variable V making the mutual information $I(X, Y; \Pi(X, Y) \mid V)$ large for any correct protocol Π , the conditional information complexity of h is necessarily large.

With this, we obtain our main direct sum theorem for information complexity.

Theorem 11.3.13. *Let μ be a fooling distribution $\mathcal{X} \times \mathcal{Y}$ for a function f with primitive h . Then*

$$\text{IC}_\delta(f) \geq n \cdot \text{CIC}_\delta^\mu(h).$$

Proof Let $V = V_1^n \in \mathcal{V}^n$ be any random vector with i.i.d. entries making (X_i, Y_i) conditionally independent given V_i . Then for any protocol Π , we have

$$\begin{aligned} I(X_1^n, Y_1^n; \Pi) &= H(\Pi) - H(\Pi \mid X_1^n, Y_1^n) \\ &= H(\Pi) - H(\Pi \mid X_1^n, Y_1^n, V) \geq H(\Pi \mid V) - H(\Pi \mid X_1^n, Y_1^n, V) = I(X_1^n, Y_1^n; \Pi \mid V) \end{aligned}$$

because we have the Markov chain $V \rightarrow (X_1^n, Y_1^n) \rightarrow \Pi$. Using the chain rule for mutual information, where we recognize that X_1^n and Y_1^n are independent given V , we have

$$\begin{aligned} I(X_1^n, Y_1^n; \Pi \mid V) &= \sum_{i=1}^n I(X_i, Y_i; \Pi \mid V, X_1^{i-1}, Y_1^{i-1}) \\ &= \sum_{i=1}^n H(X_i, Y_i \mid V, X_1^{i-1}, Y_1^{i-1}) - H(X_i, Y_i \mid V, \Pi, X_1^{i-1}, Y_1^{i-1}) \\ &\geq \sum_{i=1}^n H(X_i, Y_i \mid V) - H(X_i, Y_i \mid V, \Pi) = \sum_{i=1}^n I(X_i, Y_i; \Pi \mid V) \end{aligned} \quad (11.3.5)$$

because conditioning reduces entropy and (X_i, Y_i) are independent of X_1^{i-1}, Y_1^{i-1} given V .

Now we come to the key reduction from the global protocol Π to one solving individual primitives. On inputs $(x, y) \in \mathcal{X} \times \mathcal{Y}$, define the simulated protocol $\Pi_{i,v}(x, y)$ so that given the vector $v_{\setminus i} \in \mathcal{V}^{n-1}$, Alice and Bob independently generate $(X_j^*, Y_j^*) \stackrel{\text{iid}}{\sim} \mu(\cdot \mid V_j = v_j)$ for $j \neq i$, which is possible because of the assumed conditional independence given V , yielding $X_{\setminus i}^* \in \mathcal{X}^{n-1}$ and $Y_{\setminus i}^* \in \mathcal{Y}^{n-1}$, respectively. They then execute the protocol $\Pi((X_{\setminus i}^*, x), (Y_{\setminus i}^*, y))$ (where we substitute x and y into input position i for each). Two key consequences of this simulation follow: that $\Pi_{i,v}$ is a δ -error protocol for the primitive h and that we have the distributional equality

$$(X_i, Y_i, V_i, \Pi_{i,v}(X_i, Y_i)) \stackrel{\text{dist}}{=} (X_i, Y_i, V_i, \Pi(X_1^n, Y_1^n)) \mid V_{\setminus i} = v_{\setminus i}, \quad (11.3.6)$$

that is, the joint over the simulated protocol is equal to that over the original protocol Π conditional on $V_{\setminus i} = v_{\setminus i}$. The latter claim (11.3.6) is essentially definitional; the former requires a bit more work.

To see that $\Pi_{i,v}$ is a δ -error protocol for the primitive h , note that by construction, $X_{\setminus i}^*$ and $Y_{\setminus i}^*$ are in the support of μ , and so admit planted solutions. In particular, $f((X_{\setminus i}^*, x), (Y_{\setminus i}^*, y)) = h(x, y)$, and so $\Pi_{i,v}$ is necessarily a δ -error protocol.

The distributional equality (11.3.6) guarantees that for any v we have

$$I(X_i, Y_i; \Pi(X_1^n, Y_1^n) \mid V_i, V_{\setminus i} = v_{\setminus i}) = I(X_i, Y_i; \Pi_{i,v}(X_i, Y_i) \mid V_i),$$

and as $\Pi_{i,v}$ is a δ -error protocol for h , we have

$$\inf_v I(X_i, Y_i; \Pi_{i,v}(X_i, Y_i) \mid V_i) \geq \text{CIC}_\delta^\mu(h).$$

Substituting in the bound (11.3.5), we obtain

$$I(X_1^n, Y_1^n; \Pi) \geq \sum_{i=1}^n I(X_i, Y_i; \Pi \mid V) \geq \sum_{i=1}^n \inf_v I(X_i, Y_i; \Pi_{i,v}(X_i, Y_i) \mid V_i) \geq n \text{CIC}_\delta^\mu(h),$$

as desired. \square

With Theorem 11.3.13 in hand, we have our desired direct sum result, so that proving information complexity lower bounds reduces to providing lower bounds on the (conditional) information complexity of various 1-bit primitives. The following corollary highlights the theorem's applications to inner product and set disjointness (Examples 11.3.10 and 11.3.11).

Corollary 11.3.14. *Let f be the binary inner product $f(x, y) = \langle x, y \rangle \bmod 2$ or the disjointness function $f(x, y) = \mathbf{1}\{\langle x, y \rangle > 0\}$. Let μ be the fooling distribution in Example 11.3.12. Then*

$$\text{IC}_\delta(f) \geq n \cdot \text{CIC}_\delta^\mu(h)$$

where $h(a, b) = ab$ is the product (or AND) function.

Exercise 11.10 explores similar techniques for the entrywise lesser than or equal function, showing similar complexity lower bounds.

11.3.4 The structure of randomized communication and communication complexity of primitives

Theorem 11.3.13 provides a powerful direct sum result that demonstrates that, at least if a problem admits planted solutions for (nearly) i.i.d. sampling, then the information complexity must scale at least linearly in the complexity of the primitives making up the function f . Thus, we turn to providing information lower bounds for computing different primitive functions. Our main tool will be to show that even randomized communication protocols essentially partition the input space $\mathcal{X} \times \mathcal{Y}$ into rectangles—in analogy with Proposition 11.3.2 in the deterministic case—which allows us to provide lower bounds. The broad idea is simple: if we have an accurate protocol for computing a certain function h , we must necessarily be able to distinguish between the distribution of Π on different inputs (x, y) , as the fundamental connection between tests and variation distance (Proposition 2.3.1) reveals.

Our main goal now is to prove the following proposition, which gives a lower bound on the (conditional) information complexity of computing the AND of two bits.

Proposition 11.3.15. *Let $h(x, y) = xy$ for inputs $x, y \in \{0, 1\}$. Let μ be the fooling distribution in Example 11.3.12. Then*

$$\text{CIC}_\delta^\mu(h) \geq \frac{1}{4} \left(1 - 2\sqrt{\delta(1-\delta)}\right).$$

We prove this proposition in the remainder of this section, noting that as an immediate corollary, we obtain the following lower bounds on the communication complexity of set disjointness and binary inner product.

Corollary 11.3.16. *Let f be the binary inner product $f(x, y) = \langle x, y \rangle \bmod 2$ or the disjointness function $f(x, y) = \mathbf{1}\{\langle x, y \rangle > 0\}$. Then*

$$\text{IC}_\delta(f) \geq \frac{n}{4} (1 - 2\sqrt{\delta(1-\delta)}).$$

To control the complexity of computing individual primitives, it proves easier to use metrics tied more directly to testing. To that end, we recall the connection between Hellinger distance and the mutual information, or Jensen-Shannon divergence, between a variable X and a single bit $B \in \{0, 1\}$ in Proposition 2.2.10, which gives that if $B \rightarrow Z$, where $Z \sim P_b$ conditional on $B = b$, then

$$I_2(Z; B) \geq d_{\text{hel}}^2(P_0, P_1).$$

To apply this inequality, recall the fooling distribution μ for inner products in Example 11.3.12, where $V \sim \text{Uniform}\{0, 1\}$ and conditional on $V = 0$ we set $X = 0$ and draw $Y \sim \text{Uniform}\{0, 1\}$, and otherwise $Y = 0$ and $X \sim \text{Uniform}\{0, 1\}$. Then for $V \rightarrow (X, Y)$ from this distribution, we have

$$I_2(X, Y; \Pi(X, Y) | V) = \frac{1}{2} I_2(Y; \Pi(0, Y) | V = 0) + \frac{1}{2} I_2(X; \Pi(X, 0) | V = 1).$$

Letting Q_{xy} denote the (conditional) distribution over Π on input bits $x, y \in \{0, 1\}$ and noting that X and Y above are each uniform on $\{0, 1\}$, we see that Proposition 2.2.10 applies and so

$$I_2(X, Y; \Pi(X, Y) | V) \geq \frac{1}{2} d_{\text{hel}}^2(Q_{01}, Q_{00}) + \frac{1}{2} d_{\text{hel}}^2(Q_{10}, Q_{00}).$$

Applying the triangle inequality that $(a - b)^2 \leq (|a - c| + |c - b|)^2 \leq 2(a - c)^2 + 2(b - c)^2$, we obtain the following lemma.

Lemma 11.3.17. *Let Π be any protocol acting on two bit inputs $x, y \in \{0, 1\}$, and let μ be the fooling distribution in Example 11.3.12. Let Q_{xy} be the distribution of $\Pi(x, y)$ on inputs x, y . Then*

$$I_2(X, Y; \Pi(X, Y) | V) \geq \frac{1}{4} d_{\text{hel}}^2(Q_{01}, Q_{10}).$$

The last step in the proof of Proposition 11.3.15 is to demonstrate a property of (randomized) protocols Π analogous to the rectangular property of deterministic communication that Propositions 11.3.2 and 11.3.5 demonstrate. In analogy with the output leaf in the tree for deterministic communication complexity, let τ be the *transcript* of the communication protocol, that is, its entire communication trace. Then we claim the following analog of Proposition 11.3.2 that the set of inputs resulting in a particular output in deterministic complexity is a rectangle in $\mathcal{X} \times \mathcal{Y}$.

Lemma 11.3.18. *Let Π be any randomized protocol with inputs in $\mathcal{X} \times \mathcal{Y}$. Then there exist functions q_x and q_y such that for any transcript τ ,*

$$\mathbb{P}(\Pi(x, y) = \tau) = q_x(\tau) \cdot q_y(\tau).$$

Proof We may view any randomized protocol as a particular instantiation of a deterministic protocol $\Pi(\cdot, \cdot, u_a, u_b)$, where $u_a, u_b \in [0, 1]$ are realizations of the randomness available to Alice and Bob, respectively, inducing a particular binary communication tree. By Proposition 11.3.2, for any leaf l , the set

$$R_l(u_a, u_b) = \{(x, y) \in \mathcal{X} \times \mathcal{Y} \mid \Pi(x, y, u_a, u_b) \text{ reaches } l\}$$

is a rectangle, that is, $R_l(u_a, u_b) = A_l(u_a) \times B_l(u_b)$ for sets $A_l(u) \subset \mathcal{X}$ and $B_l(u) \subset \mathcal{Y}$. Of course, the leaves l of the tree are in bijection with the entire transcript τ , so that if τ ends in leaf l , then

$$\mathbb{P}(\Pi(x, y) = \tau) = \mathbb{P}((x, y) \in R_l(U_a, U_b)) = \mathbb{P}(x \in A_l(U_a), y \in B_l(U_b))$$

where $U_a, U_b \stackrel{\text{iid}}{\sim} \text{Uniform}[0, 1]$ are the the randomness Alice and Bob use, respectively.

Expanding this as an integral gives

$$\begin{aligned} \mathbb{P}(x \in A_l(U_a), y \in B_l(U_b)) &= \int_0^1 \int_0^1 \mathbf{1}\{x \in A_l(u_a)\} \mathbf{1}\{y \in B_l(u_b)\} du_a du_b \\ &= \mathbb{P}(x \in A_l(U_a)) \mathbb{P}(y \in B_l(U_b)). \end{aligned}$$

Set $q_x(\tau) = \mathbb{P}(x \in A_l(U_a))$ and $q_y(\tau) = \mathbb{P}(y \in B_l(U_b))$. □

We thus have the following key *cut and paste* property, which shows that in some sense, Hellinger distances respect the “rectangular” structure of communication protocols.

Lemma 11.3.19. *Let Π be any protocol acting on inputs in $\mathcal{X} \times \mathcal{Y}$ and let $Q_{x,y}$ be the distribution of $\Pi(x, y)$ on inputs x, y . Then*

$$d_{\text{hel}}(Q_{x,y}, Q_{x',y'}) = d_{\text{hel}}(Q_{x,y'}, Q_{x',y}).$$

Proof Let \mathcal{T} be the collection of all possible transcripts the protocol outputs. By Lemma 11.3.18 we have

$$\begin{aligned} d_{\text{hel}}^2(Q_{x,y}, Q_{x',y'}) &= \frac{1}{2} \sum_{\tau \in \mathcal{T}} \left(\sqrt{Q_{x,y}(\tau)} - \sqrt{Q_{x',y'}(\tau)} \right)^2 \\ &= \frac{1}{2} \sum_{\tau \in \mathcal{T}} \left(\sqrt{q_x(\tau)q_y(\tau)} - \sqrt{q_{x'}(\tau)q_{y'}(\tau)} \right)^2 = 1 - \sum_{\tau} \sqrt{q_x(\tau)q_y(\tau)q_{x'}(\tau)q_{y'}(\tau)}. \end{aligned}$$

Rearranging by the trivial modification $q_x q_y q_{x'} q_{y'} = q_x q_{y'} q_{x'} q_y$, we have the result. □

We now finalize the proof of Proposition 11.3.15. Substituting this cutting and pasting in Lemma 11.3.17 we have

$$I_2(X, Y; \Pi(X, Y) \mid V) \geq \frac{1}{4} d_{\text{hel}}^2(Q_{01}, Q_{10}) = \frac{1}{4} d_{\text{hel}}^2(Q_{00}, Q_{11}).$$

Then a simple lemma recalling the testing inequalities in Chapter 2.3.1 completes the proof of the proposition, because it guarantees that $4I_2(X, Y; \Pi(X, Y) \mid V) \geq 1 - 2\sqrt{\delta(1-\delta)}$ no matter the choice of protocol Π , and so

$$\text{CIC}_{\delta}^{\mu}(h) \geq \inf_{\Pi} I_2(X, Y; \Pi(X, Y) \mid V) \geq \frac{1}{4} \left(1 - 2\sqrt{\delta(1-\delta)} \right).$$

Lemma 11.3.20. *Let Π be any δ -accurate protocol for computing $h(x, y) = xy$ and Q_{xy} be its distribution on inputs (x, y) . Then $d_{\text{hel}}^2(Q_{00}, Q_{11}) \geq 1 - 2\sqrt{\delta(1 - \delta)}$.*

Proof Assume that Π computes the product $xy \in \{0, 1\}$ correctly with probability at least $1 - \delta$, that is, $\mathbb{P}(\Pi_{\text{out}}(x, y) \neq xy) \leq \delta$ for all $x, y \in \{0, 1\}$. By Le Cam's testing lower bounds (Proposition 2.3.1), we know that

$$\begin{aligned} 2\delta &\geq \mathbb{P}(\Pi_{\text{out}}(0, 0) \neq 0) + \mathbb{P}(\Pi_{\text{out}}(1, 1) \neq 1) \geq 1 - \|Q_{00} - Q_{11}\|_{\text{TV}} \\ &\stackrel{(*)}{\geq} 1 - d_{\text{hel}}(Q_{00}, Q_{11}) \sqrt{2 - d_{\text{hel}}^2(Q_{00}, Q_{11})}, \end{aligned}$$

where inequality $(*)$ follows from the inequalities in Proposition 2.2.7 relating Hellinger and total-variation distance. Let $d = d_{\text{hel}}^2(Q_{00}, Q_{11})$ for shorthand. Then rearranging gives $d(2 - d) \geq (1 - 2\delta)^2$. Solving for d in $0 \geq d^2 - 2d + (1 - 2\delta)^2$ yields $d \geq 1 - \sqrt{1 - (1 - 2\delta)^2}$. Recognize that $1 - (1 - 2\delta)^2 = 4(\delta - \delta^2)$. \square

11.4 Communication complexity in estimation

A major application combining strong data processing inequalities and communication is in the communication and information complexity of statistical estimation itself. In this context, we limit the amount of information—or perhaps bits—that a procedure may send about individual examples, and then ask to what extent this constrains the estimator. This has applications in situations in which the memory available to an estimator is limited, in situations with privacy—as we shall see—and of course, when we restrict the number of bits different machines storing distributed data may send.

We consider the following setting: m machines, or agents, have data X_i , $i = 1, \dots, m$. Communication proceeds in rounds $t = 1, 2, \dots, T$, where in each round t machine i sends datum $Z_i^{(t)}$. To allow for powerful protocols—with little restriction except that each machine i may send only a certain amount of information—we allow $Z_i^{(t)}$ to depend arbitrarily on the previous messages $Z_1^{(t)}, \dots, Z_{i-1}^{(t)}$ as well as $Z_k^{(\tau)}$ for all $k \in \{1, \dots, m\}$ and $\tau < t$. We visualize this as a public blackboard B , where in each round t each $Z_i^{(t)}$ is collected into $B^{(t)}$, along with the previous public blackboards $B^{(\tau)}$ for $\tau < t$, and all machines may read these public blackboards. Thus, in round t , individual i generates the communicated variable $Z_i^{(t)}$ according to the channel

$$Q_{Z_i^{(t)}}(\cdot \mid X_i, Z_{<i}^{(t)}, B^{(t-1)}) = Q_{Z_i^{(t)}}(\cdot \mid X_i, Z_{\rightarrow i}^{(t)}).$$

Here we have used the notation $Z_{<i} := (Z_1, \dots, Z_{i-1})$, and we will use $Z_{\leq i} := (Z_1, \dots, Z_i)$ and similarly for superscripts throughout. We will also use the notation $Z_{\rightarrow i}^{(t)} = (B^{(1)}, Z_{<i}^{(t)})$ to denote all the messages coming into communication of $Z_i^{(t)}$. Figure 11.3 illustrates two rounds of this communication scheme.

We can provide lower bounds on the minimax risk of communication-constrained estimators by extending the data processing inequality approach we have developed. Our approach to the lower bounds, which we provide in Sections 11.4.1 and 11.4.2 to follow, is roughly as follows. First, we develop another *direct sum* bound, in analogy with Theorem 11.3.13, meaning that the difficulty of

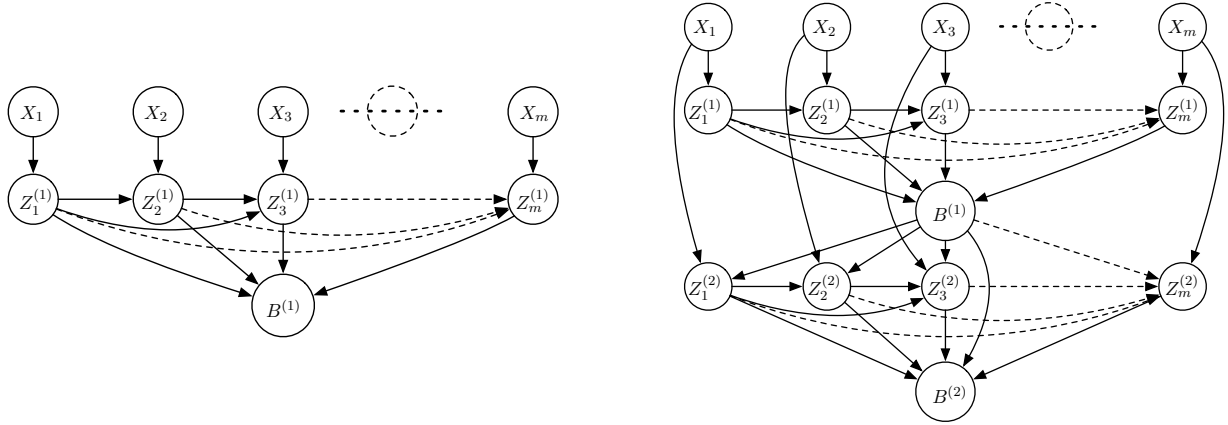


Figure 11.3. Left: single round of communication of variables, writing to public blackboard $B^{(1)}$. Right: two rounds of communication of variables, writing to public blackboards $B^{(1)}$ and $B^{(2)}$.

solving a d -dimensional problem is roughly d -times that of solving a 1-dimensional version of the problem; thus, any lower bounds on the error in 1-dimensional problems imply lower bounds for d -dimensional problems. Second, we provide an extension of the data processing inequalities we have developed thus far to apply to particular communication scenarios.

The key to our reductions is that we consider families of distributions where the coordinates of X are independent, which dovetails with Assouad's method. We thus index our distributions by $v \in \{0, 1\}^d$, and in proving our lower bounds, we assume the typical Markov structure

$$V \rightarrow (X_1, \dots, X_m) \rightarrow \Pi(X_1^m),$$

where V is chosen uniformly at random from $\{-1, 1\}^d$, and $\Pi = \Pi(X_1^m)$ denotes the protocol of the entire communication—in this context, this is the entire set of blackboard messages

$$\Pi = (B^{(1)}, \dots, B^{(T)}),$$

(which also encodes the message order). We assume that X follows a d -dimensional product distribution, so that conditional on $V = v$ we have

$$X \stackrel{\text{iid}}{\sim} P_v = P_{v_1} \otimes P_{v_2} \otimes \dots \otimes P_{v_d}. \quad (11.4.1)$$

The generation strategy (11.4.1) guarantees that conditional on the j th coordinate $V_j = v_j$, the coordinates $X_{i,j}$ are i.i.d. and independent of $V_{\setminus j} = (V_1, \dots, V_{j-1}, V_{j+1}, \dots, V_d)$ as well as independent of $X_{i',j}$ for data points $i' \neq i$.

11.4.1 Direct sum communication bounds

Our first step is to argue that, if we can prove a lower bound on the information complexity of one-dimensional estimation, we can prove a lower bound on d -dimensional problems that scales with the dimension. To accomplish this reduction, let $X_{\leq m,j} = (X_{i,j})_{i=1}^m$ be the j th coordinate of the data, and let $X_{\leq m,\setminus j}$ be the remaining $d-1$ coordinates across all $i = 1, \dots, m$. Then by the construction (11.4.1), we have the Markov structure

$$V_j \rightarrow X_{\leq m,j} \rightarrow \Pi(X_1^m) \leftarrow X_{\leq m,\setminus j} \leftarrow V_{\setminus j}.$$

In particular, viewing $X_{\leq m, \setminus j}$ as extraneous randomness, we have the simpler Markovian structure

$$V_j \rightarrow X_{\leq m, j} \rightarrow \Pi, \quad (11.4.2)$$

so that we may think of the communication $\Pi = \Pi(X_{\leq m, j})$ as acting only on $X_{\leq m, j}$. Now, define M_{-j} and M_j to be the marginal distributions over the total communication protocol Π conditional on $V_j = \pm j$, the one-variable model (11.4.2). Then Le Cam's testing equality (Proposition 2.3.1), and the equivalence between Hellinger and variation distance (Proposition 2.2.7) imply that

$$\begin{aligned} \inf_{\hat{V}} 2 \sum_{j=1}^d \mathbb{P}(\hat{V}_j(\Pi) \neq V_j) &\geq \sum_{j=1}^d (1 - \|M_{-j} - M_{+j}\|_{TV}) \geq \sum_{j=1}^d (1 - \sqrt{2} d_{\text{hel}}(M_{-j}, M_{+j})) \\ &\geq d \left(1 - \sqrt{\frac{2}{d} \sum_{j=1}^d d_{\text{hel}}^2(M_{-j}, M_{+j})} \right) \end{aligned}$$

by Cauchy-Schwarz. Summarizing, we have the following

Proposition 11.4.1 (Assouad's method in communication). *Let M_{+j} be the marginal distribution over Π conditional on $V_j = 1$ and M_{-j} be the marginal distribution of Π conditional on $V_j = -1$ in Markov structure (11.4.2) and assume X_i follow the product distribution (11.4.1). Then*

$$\sum_{j=1}^d \mathbb{P}(\hat{V}_j(\tau) \neq V_j) \geq \frac{d}{2} \left(1 - \sqrt{\frac{2}{d} \sum_{j=1}^d d_{\text{hel}}^2(M_{-j}, M_{+j})} \right).$$

Recalling Assouad's method (Lemma 9.5.2) of Chapter 9.5, we see that any time we have a problem with separation with respect to the Hamming metric (9.5.1), we have a lower bound on its error in estimation problems. This proposition analogizes Theorem 11.3.13, in that small Hellinger distance between the individual marginals $M_{\pm j}$ necessarily makes the testing and estimation problems hard.

11.4.2 Communication data processing

We now revisit the data processing inequalities in Section 11.1, where we consider a variant that allows us to prove lower bounds for estimation problems with limited communication. It will be more notationally convenient in this section to use $V \in \{0, 1\}$ rather than $\{-1, 1\}$, so we do so without comment. Our starting point is a revised strong data processing inequality.

Definition 11.7. *Let P_0, P_1 be arbitrary distributions on a space \mathcal{X} , let $V \in \{0, 1\}$ uniformly at random, and conditional on $V = v$, draw $X \sim P_v$. Consider the Markov chain $V \rightarrow X \rightarrow Z$. The mutual information strong data processing constant $\beta(P_0, P_1)$ is*

$$\beta(P_0, P_1) := \sup_{X \rightarrow Z} \frac{I(V; Z)}{I(X; Z)},$$

where the supremum is taken over all conditional distributions (Markov kernels) from X to Z .

In contrast to Definition 11.1, in this definition we have a contraction over the “beginning” of the chain $V \rightarrow X$ rather than the distribution $X \rightarrow Z$. Identifying Z with a communication protocol $\Pi(X_1^m)$, this makes it possible to develop lower bounds on estimation and testing that then depend on the information $I(X; \Pi)$.

Distributions with bounded likelihood ratios provide one way to demonstrate a strong data processing inequality of the form in Definition 11.7, where in analogy with Theorem 11.2.1 we obtain a contraction inequality involving the total variation distance.

Proposition 11.4.2. *Let $V \rightarrow X \rightarrow Z$, where $X \sim P_v$ conditional on $V = v$. Let P_X and $P_X(\cdot | Z)$ denote the marginal and conditional distributions on X given Z , respectively. If $|\log \frac{dP_v}{dP_{v'}}| \leq \alpha$ for all v, v' , then*

$$I(V; Z) \leq 4(e^\alpha - 1)^2 \mathbb{E}_Z \left[\|P_X(\cdot | Z) - P_X\|_{\text{TV}}^2 \right] \leq 2(e^\alpha - 1)^2 I(X; Z).$$

We leave the proof of this proposition as Exercise 11.12, as it follows by adapting the techniques we use to prove Theorem 11.2.1, with the main difference being the random variables with bounded likelihood ratios ($X \rightarrow Z$ versus $V \rightarrow X$). A brief example illustrates Proposition 11.4.2.

Example 11.4.3 (Bernoulli distributions): Let $P_v = \text{Bernoulli}(\frac{1+v\delta}{2})$ for $v \in \{-1, 1\}$. Then we have likelihood ratio bound

$$\left| \log \frac{dP_1}{dP_{-1}} \right| \leq \log \frac{1+\delta}{1-\delta}$$

and so under the conditions of Proposition 11.4.2, for any Z we have

$$I(V; Z) \leq 2 \left(\frac{1+\delta}{1-\delta} - 1 \right)^2 I(X; Z) = 2 \left(\frac{2\delta}{1-\delta} \right)^2 I(X; Z) \stackrel{(i)}{\leq} 10\delta^2 I(X; Z),$$

where inequality (i) holds for $\delta \in [0, 1/10]$. \diamond

We now give the two main results connecting mutual information and the contraction-type bounds in Definition 11.7. To provide bounds using Proposition 11.4.1, we wish to control the Hellinger distance between individual marginals $M_{\pm j}$, so we consider single variables in the Markov chain

$$V \rightarrow (X_1, \dots, X_m) \rightarrow \Pi,$$

where $V \in \{0, 1\}$. To state the coming theorems, we make a restriction on the data generation $V \rightarrow X$, calling distributions P_0 and P_1 (c, β) -contractive if

$$\beta(P_0, P_1) \leq \beta \leq 1 \quad \text{and} \quad \max \{D_\infty(P_0 \| P_1), D_\infty(P_1 \| P_0)\} \leq \log c, \quad (11.4.3)$$

where $D_\infty(\cdot \| \cdot)$ denotes the Rényi- ∞ -divergence. Proposition 11.4.2 shows that whenever such a c exists we certainly have $\beta(P_0, P_1) \leq 2(c-1)^2$.

The next theorem then provides the basic information contraction inequality for single-variable communication.

Theorem 11.4.4. *Let $1 \leq c < \infty$ and $\beta \leq 1$. Let P_0 and P_1 be (c, β) -contractive (11.4.3) distributions on \mathcal{X} and M_v , $v \in \{0, 1\}$ be the marginal distribution of the protocol Π conditional on $V = v$. Then*

$$d_{\text{hel}}^2(M_0, M_1) \leq \frac{7}{2}(c+1)\beta \cdot \min \{I(X_1^m; \Pi(X_1^m) | V = 0), I(X_1^m; \Pi(X_1^m) | V = 1)\}.$$

The proof of Theorem 11.4.4 is quite complicated, so we defer it to Section 11.5.

We can use Theorem 11.4.4 to obtain bounds on the probability of error—detection of d -dimensional signals—in higher dimensional problems based on mutual information alone. Because the theorem provides a bound involving the minimum of the conditional mutual informations, we have substantial freedom to combine the direct-sum lower bounds in Section 11.4.1 to massage it into the mutual information between the data X_1^m and the protocol $\Pi(X_1^m)$.

We thus recall the definition (11.4.1) of our product distribution signals, where we assume that each individual datum $X_i = (X_{i,1}, \dots, X_{i,d}) = (X_{i,j})_{j=1}^d$ belongs to a d -dimensional set and conditional on $V = v \in \{-1, 1\}^d$ has independent coordinates distributed as $X_{i,j} \sim P_{v_j}$. With this, we have the following theorem, which follows by a combination of Assouad's method (in the context of communication bounds, i.e. Proposition 11.4.1) and Theorem 11.4.4.

Theorem 11.4.5. *Let Π the entire communication protocol in Figure 11.3, $V \in \{-1, 1\}^d$ be uniform, and generate $X_i \stackrel{\text{iid}}{\sim} P_v$, $i = 1, \dots, m$ according to the independent coordinate distribution (11.4.1). Assume additionally that for each coordinate $j = 1, \dots, d$, the coordinate distributions $P_{\pm v_j}$ are (c, β) -contractive (11.4.3). Then for any estimator \hat{V} ,*

$$\sum_{j=1}^d \mathbb{P}(\hat{V}_j(\Pi) \neq V_j) \geq \frac{d}{2} \left(1 - \sqrt{7(c+1) \frac{\beta}{d} \cdot I(X_1, \dots, X_m; \Pi \mid V)} \right).$$

Proof Under the given conditions, Proposition 11.4.1 and Theorem 11.4.4 immediately combine to give

$$\sum_{j=1}^d \mathbb{P}(\hat{V}_j(\Pi) \neq V_j) \geq \frac{d}{2} \left(1 - \sqrt{7(c+1) \frac{\beta}{d} \sum_{j=1}^d \min_{v \in \{-1, 1\}} I(X_{1,j}, \dots, X_{m,j}; \Pi \mid V_j = v)} \right).$$

Certainly

$$\min_{v \in \{-1, 1\}} I(X_{1,j}, \dots, X_{m,j}; \Pi \mid V_j = v) \leq I(X_{1,j}, \dots, X_{m,j}; \Pi \mid V_j).$$

Then, using that w.l.o.g. we may assume the $X_{i,j}$ are discrete, we obtain

$$\begin{aligned} \sum_{j=1}^d I((X_{i,j})_{i=1}^m; \Pi \mid V_j) &= \sum_{j=1}^d [H((X_{i,j})_{i=1}^m \mid V_j) - H((X_{i,j})_{i=1}^m \mid \Pi, V_j)] \\ &\stackrel{(i)}{=} \sum_{j=1}^d [H((X_{i,j})_{i=1}^m \mid (X_{i,j'})_{i \leq m, j' < j}, V) - H((X_{i,j})_{i=1}^m \mid \Pi, V_j)] \\ &\leq \sum_{j=1}^d [H((X_{i,j})_{i=1}^m \mid (X_{i,j'})_{i \leq m, j' < j}, V) - H((X_{i,j})_{i=1}^m \mid (X_{i,j'})_{i \leq m, j' < j}, \Pi, V)] \\ &= \sum_{j=1}^d I((X_{i,j})_{i=1}^m; \Pi \mid V, (X_{i,j'})_{i \leq m, j' < j}) = I(X_1, \dots, X_m; \Pi \mid V), \end{aligned}$$

where equality (i) used the independence of $X_{i,j}$ from $V_{\setminus j}$ and $X_{i,j'}$ for $j' \neq j$ given V_j , and the inequality that conditioning reduces entropy. This gives the theorem. \square

11.4.3 Applications: communication and privacy lower bounds

Let us now turn to a few different applications of our lower bounds on communication-constrained estimators. We evidently require two conditions: first, we must show that the distributions our data follows satisfy a strong (mutual information) data processing inequality (Definition 11.7). Second,

we must provide a (good enough) upper bound on the mutual information $I(X_1, \dots, X_m; \Pi \mid V)$ between the data points X_i and communication protocol. While there are many strategies to providing bounds and strong data processing inequalities, we focus mainly on situations with bounded likelihood ratio, where Proposition 11.4.2 directly provides the type of strong data processing inequality we require.

Communication lower bounds

Our first set of examples considers direct communication bounds, where controlling $I(X_1^m; \Pi)$ is relatively straightforward. Assume the setting in the introduction to Section 11.4, where to establish our communication bounds we assume each machine $i = 1, \dots, m$ may send at most B_i total bits of information throughout the entire communication protocol—that is, for each pair i, t , we have a bound

$$H(Z_i^{(t)} \mid Z_{\rightarrow i}^{(t)}) \leq B_{i,t} \quad \text{and} \quad \sum_t B_{i,t} \leq B_i \quad (11.4.4)$$

on the message from X_i in round t . (This is a weaker condition than $H(Z_i^{(t)}) \leq B_{i,t}$ for each i, t .) With this bound, we can provide minimax lower bounds on communication-constrained estimator.

For our first collection, we consider estimating the parameters of d independent Bernoulli distributions in squared error. Let \mathcal{P}_d be the family of d -dimensional Bernoulli distributions, where we let the parameter $\theta \in [0, 1]^d$ be such that $P_\theta(X_j = 1) = \theta_j$. Then we have the following result.

Proposition 11.4.6. *Let $\mathfrak{M}_m(\theta(\mathcal{P}_d), \|\cdot\|_2^2, \{B_i\}_{i=1}^m)$ denote the minimax mean-square error for estimation of a d -dimensional Bernoulli under the information constraint (11.4.4). Then*

$$\mathfrak{M}_m(\theta(\mathcal{P}_d), \|\cdot\|_2^2, \{B_i\}_{i=1}^m) \geq c \min \left\{ \frac{d}{m} \frac{d}{\frac{1}{m} \sum_{i=1}^m B_i}, d \right\},$$

where $c > 0$ is a numerical constant.

Proof By the standard Assouad reduction (Section 9.5), when we take coordinate distributions $P_{v_j} = \text{Bernoulli}(\frac{1+\delta v_j}{2})$, we have a $c\delta^2$ -separation in Hamming metric. Applying Theorem 11.4.5 and Example 11.4.3, we obtain the minimax lower bound, valid for $0 \leq \delta \leq \frac{1}{10}$, of

$$\mathfrak{M}_m(\theta(\mathcal{P}_d), \|\cdot\|_2^2, \{B_i\}_{i=1}^m) \geq c\delta^2 d \left(1 - \sqrt{C \frac{\delta^2}{d} I(X_1, \dots, X_m; \Pi \mid V)} \right).$$

Now, we note that for any Markov chain $V \rightarrow X \rightarrow Z$,

$$I(X; Z \mid V) = H(Z \mid V) - H(Z \mid X, V) = H(Z \mid V) - H(Z \mid X) \leq H(Z) - H(Z \mid X) = I(X; Z).$$

Thus we obtain

$$\begin{aligned} I(X_1, \dots, X_m; \Pi \mid V) &\leq I(X_1, \dots, X_m; \Pi) \\ &= \sum_{i=1}^m \sum_{t=1}^T I(X_1, \dots, X_m; Z_i^{(t)} \mid Z_{\rightarrow i}^{(t)}). \end{aligned}$$

As the message $Z_i^{(t)}$ satisfies the conditional independence $Z_i^{(t)} \perp\!\!\!\perp X_{\setminus i} \mid Z_{\rightarrow i}^{(t)}, X_i$, this final quantity equals $\sum_{i,t} I(X_i; Z_i^{(t)} \mid Z_{\rightarrow i}^{(t)})$. But of course $I(X_i; Z_i^{(t)} \mid Z_{\rightarrow i}^{(t)}) \leq H(Z_i^{(t)} \mid Z_{\rightarrow i}^{(t)}) \leq B_{i,t}$, and so

$$\mathfrak{M}_m(\theta(\mathcal{P}_d), \|\cdot\|_2^2, \{B_i\}_{i=1}^m) \geq c\delta^2 d \left(1 - \sqrt{C \frac{\delta^2}{d} \sum_{i,t} B_{i,t}} \right).$$

Choosing $\delta = \min\{1/10, \frac{d}{2C \sum_i B_i}\}$ gives the result. \square

This result deserves some discussion. It is sharp in the case that the number of bits is of order d or less from each machine: when we set $B_i = d$, the lower bound becomes

$$\sup_{\theta} \mathbb{E}_{\theta}[\|\hat{\theta}(\Pi) - \theta\|_2^2] \gtrsim \min\left\{\frac{d}{m} \cdot \frac{d}{d}, d\right\} = \frac{d}{m},$$

which is certainly achievable (each machine simply sends its entire vector $X_i \in \{0,1\}^d$). When machines communicate fewer than d bits, we have a tighter result; for example, if only k/m machines send d bits, and the rest communicate little, we obtain

$$\sup_{\theta} \mathbb{E}_{\theta}[\|\hat{\theta}(\Pi) - \theta\|_2^2] \gtrsim \min\left\{\frac{d}{m} \cdot \frac{md}{kd}, d\right\} = \frac{d}{k},$$

which is similarly intuitive. The extension of these ideas to the case when each machine has an individual sample of size n is more challenging, as it requires tensorized variants of the strong data processing inequality in Definition 11.7; we provide remarks in the bibliographical section.

Lower bounds in locally private estimation

We return to the local privacy setting we consider in Section 11.2, except now we allow substantially more interaction. We treat local differential privacy in the communication model of Figure 11.3, where n individuals have data X_i which they wish to privatize, and proceed in rounds, releasing data $Z_i^{(t)}$ from individual i in round t . A natural setting is to assume each data release $Z_i^{(t)}$ is $\varepsilon_{i,t}$ -differentially private: instead of the sequentially interactive model (11.2.1), we have

$$Q(Z_i^{(t)} \in A \mid X_i = x, Z_{\rightarrow i}^{(t)} = z_{\rightarrow i}^{(t)}) \leq \exp(\varepsilon_{i,t}) \cdot Q(Z_i^{(t)} \in A \mid X_i = x', Z_{\rightarrow i}^{(t)} = z_{\rightarrow i}^{(t)}) \quad (11.4.5)$$

for each i, t and all possible $x, x', z_{\rightarrow i}^{(t)}$. At a more abstract level, rather than a particular privacy guarantee on each individual data release $Z_i^{(t)}$, we can assume a more global stability guarantee akin to the (average) KL-stability in interactive data analysis (Definition 6.1). Thus, let $\Pi(x_1^n)$ be the entire collection of communicated information in the protocol in Figure 11.3 on input data x_1, \dots, x_n . Abusing notation to let $D_{\text{kl}}(Z_0 \parallel Z_1)$ be the KL-divergence between the distributions of Z_0 and Z_1 , as in Definition 6.1, we make the following definition to capture arbitrary interactions.

Definition 11.8 (Average KL-privacy). *Let the samples $x_{\leq n} \in \mathcal{X}^n$ and $x_{\leq n}^{(i)} \in \mathcal{X}^n$ differ only in example i . Then the data release Π is ε_{kl} -KL-locally-private on average if*

$$\frac{1}{n} \sum_{i=1}^n D_{\text{kl}}\left(\Pi(x_{\leq n}^{(i)}) \parallel \Pi(x_{\leq n})\right) \leq \varepsilon_{\text{kl}}.$$

The following observation shows that for appropriate choices of ε_{kl} , this is indeed weaker than the interactive guarantee (11.4.5).

Lemma 11.4.7. *Let the communication Q satisfy the interactive privacy guarantee (11.4.5) and Π be the induced communication protocol over rounds $t \leq T$. Then*

$$\frac{1}{n} \sum_{i=1}^n D_{\text{kl}} \left(\Pi(x_{\leq n}^{(i)}) \| \Pi(x_{\leq n}) \right) \leq \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^T \min \left\{ \varepsilon_{i,t}, \frac{3}{2} \varepsilon_{i,t}^2 \right\}.$$

Proof Using the chain rule for the KL-divergence, we have for any j that

$$\begin{aligned} D_{\text{kl}} \left(\Pi(x_{\leq n}^{(j)}) \| \Pi(x_{\leq n}) \right) &= \sum_{i=1}^n \sum_{t=1}^T \mathbb{E} \left[D_{\text{kl}} \left(Q(Z_i^t \in \cdot \mid x_i^{(j)}, Z_{\rightarrow i}^{(t)}) \| Q(Z_i^t \in \cdot \mid x_i, Z_{\rightarrow i}^{(t)}) \right) \right] \\ &= \sum_{t=1}^T \mathbb{E} \left[D_{\text{kl}} \left(Q(Z_i^t \in \cdot \mid x_j^{(j)}, Z_{\rightarrow i}^{(t)}) \| Q(Z_i^t \in \cdot \mid x_j, Z_{\rightarrow i}^{(t)}) \right) \right], \end{aligned}$$

where the expectation is taken over $Z_{\rightarrow i}^{(t)}$ in the protocol $\Pi(x_{\leq n}^{(j)})$, and the second equality follows because $x_j^{(i)} = x_j$ for all j except index i . Now let P_0 and P_1 be arbitrary distributions whose densities satisfy $p_0(z)/p_1(z) \leq e^\varepsilon$. Then

$$D_{\text{kl}}(P_0 \| P_1) \leq \varepsilon \quad \text{and} \quad D_{\text{kl}}(P_0 \| P_1) \leq \log(1 + D_{\chi^2}(P_0 \| P_1)) \leq \log(1 + (e^\varepsilon - 1)^2)$$

by Proposition 2.2.9. Then by inspection $\min\{\varepsilon, \log(1 + (e^\varepsilon - 1)^2)\} \leq \min\{\varepsilon, \frac{3}{2}\varepsilon^2\}$ for all $\varepsilon \geq 0$. Returning to the initial KL-divergence sum, we thus obtain

$$\sum_{i=1}^n D_{\text{kl}} \left(\Pi(x_{\leq n}^{(i)}) \| \Pi(x_{\leq n}) \right) \leq \sum_{i=1}^n \sum_{t=1}^T \mathbb{E} \left[\min \left\{ \varepsilon_{i,t}, \frac{3}{2} \varepsilon_{i,t}^2 \right\} \right],$$

as desired. □

The key is that the average KL-local privacy guarantee is sufficient to provide a mutual information bound, thus allowing us to apply Theorem 11.4.5 as in the proof of Proposition 11.4.6.

Proposition 11.4.8. *Let Π be any ε_{kl} -KL-locally-private on average protocol and assume that X_1, \dots, X_n are independent conditional on V . Then*

$$I(X_1, \dots, X_n; \Pi(X_1^n) \mid V) \leq n\varepsilon_{\text{kl}}.$$

Proof The conditional independence of the X_i guarantees that

$$\begin{aligned} I(X_1^n; \Pi(X_1^n) \mid V) &= \sum_{i=1}^n H(X_i \mid X_1^{i-1}, V) - H(X_i \mid \Pi, X_1^{i-1}, V) \\ &\leq \sum_{i=1}^n H(X_i \mid X_{\setminus i}, V) - H(X_i \mid \Pi, X_{\setminus i}, V) = \sum_{i=1}^n I(X_i; \Pi(X_1^n) \mid V, X_{\setminus i}). \end{aligned}$$

We abuse notation to let $\Pi^*(X_{\setminus i})$ be the marginal protocol (marginalizing over X_i). Then

$$I(X_i; \Pi(X_1^n) \mid V, X_{\setminus i}) = \mathbb{E} [D_{\text{kl}}(\Pi(X_{\setminus i}, X_i) \parallel \Pi^*(X_{\setminus i}))] \leq \mathbb{E} [D_{\text{kl}}(\Pi(X_{\setminus i}, X_i) \parallel \Pi(X_{\setminus i}, X'_i))]$$

where the first expectation is taken over V and $X_j \stackrel{\text{iid}}{\sim} P_v$ conditional on $V = v$ and the inequality uses convexity and draws X'_i independently. Summing over $i = 1, \dots, n$, Definition 11.8 gives the result. \square

Applying Theorem 11.4.5, we then obtain the following corollary.

Corollary 11.4.9. *Let the conditions of Theorem 11.4.5 hold. If the data release Π is ε_{kl} -private on average, then*

$$\sum_{j=1}^d \mathbb{P}(\hat{V}_j(\Pi) \neq V_j) \geq \frac{d}{2} \left(1 - \sqrt{7(c+1) \frac{\beta}{d} n \varepsilon_{\text{kl}}} \right).$$

Specializing to the case that we wish to estimate a d -dimensional Bernoulli vector, where $X \in \{\pm 1\}$ has coordinates with $\mathbb{P}(X_j = 1) = \theta_j$, Example 11.4.3 gives the following minimax lower bound.

Corollary 11.4.10. *Let $\mathfrak{M}_n(\theta(\mathcal{P}_d), \|\cdot\|_2^2, \varepsilon_{\text{kl}})$ denote the minimax mean-square error for estimation of a d -dimensional Bernoulli under the ε_{kl} -KL-locally-private-on-average constraint in Definition 11.8. Then*

$$\mathfrak{M}_n(\theta(\mathcal{P}_d), \|\cdot\|_2^2, \varepsilon_{\text{kl}}) \geq c \min \left\{ d, \frac{d^2}{n \varepsilon_{\text{kl}}} \right\}.$$

Proof By Corollary 11.4.9 and Example 11.4.3, we have minimax lower bound

$$\mathfrak{M}_n(\theta(\mathcal{P}_d), \|\cdot\|_2^2, \varepsilon_{\text{kl}}) \gtrsim d \delta^2 \left(1 - \sqrt{C \frac{\delta^2}{d} n \varepsilon_{\text{kl}}} \right)$$

for a numerical constant C , which is valid for $\delta \lesssim 1$. Choose δ^2 to scale as $\min\{1, \frac{d}{n \varepsilon_{\text{kl}}}\}$. \square

When instead of the average KL-privacy we use the pure local differential privacy constraint (11.4.5), Lemma 11.4.7 implies the following.

Corollary 11.4.11. *Let $\mathfrak{M}_n(\theta(\mathcal{P}_d), \|\cdot\|_2^2, \varepsilon)$ denote the minimax mean-square error for estimation of a d -dimensional Bernoulli where each data release is $\varepsilon_{i,t}$ -locally differentially private (11.4.5), and $\sum_{t=1}^{\infty} \varepsilon_{i,t} \leq \varepsilon$. Then*

$$\mathfrak{M}_n(\theta(\mathcal{P}_d), \|\cdot\|_2^2, \varepsilon) \geq c \min \left\{ d, \frac{d^2}{n(\varepsilon \wedge \varepsilon^2)} \right\}.$$

11.5 Proof of Theorem 11.4.4

The proof proceeds in stages. The basic ideas are as follows:

1. Relate the Hellinger distance between the marginal distributions M_0 and M_1 of Π conditional on $V = 0$ or 1 to a sum of Hellinger distances between the marginal M_0 and an alternative M'_i where $X_i \sim P_1$ and $X_{\setminus i} \stackrel{\text{iid}}{\sim} P_0$.
2. Provide a data processing inequality to relate $d_{\text{hel}}(M_0, M'_i)$ and the mutual information $I(X_i; \Pi)$ between the individual observation X_i and the protocol Π .
3. Use the standard chain rules for mutual information to finalize the theorem.

Step 1: sequential modification of marginals

We begin by relating the marginal distributions M_0 and M_1 by a sequence of one-variable changes. To that end, for bit vectors $b \in \{0, 1\}^m$ define M_b to be the marginal distribution over the protocol $\Pi(X_1^m)$ generated from (X_1, \dots, X_m) , where for each i we generate X_i by independently sampling

$$X_i \mid b \sim P_{b_i}. \quad (11.5.1)$$

For the standard basis vectors e_1, \dots, e_m , we expect M_0 to be close to M_{e_l} , and thus hope for some type of tensorization behavior, where we can relate M_0 and M_1 via one-step changes from M_0 to M_{e_l} . The next lemma realizes this promise.

Lemma 11.5.1. *Let M_0, M_1 , and M_{e_l} be as above. Then*

$$d_{\text{hel}}^2(M_0, M_1) \leq 7 \sum_{l=1}^m d_{\text{hel}}^2(M_0, M_{e_l}). \quad (11.5.2)$$

Proof The proof crucially relies on the Euclidean structures that the Hellinger distance induces along with analogues of the cut-and-paste (the “rectangular” structure of inputs in communication protocols) properties from deterministic and randomized two-player communication. We assume without loss of generality that Π is discrete, as the Hellinger distance is an f -divergence and so can be arbitrarily approximated by discrete random variables.

First, we analogize the “rectangular” probabilistic structure of two-player communication protocols in Lemmas 11.3.18 and 11.3.19, which yields a multi-player cut-and-paste lemma.

Lemma 11.5.2 (cutting and pasting). *Let $a, b, c, d \in \{0, 1\}^m$ be bit vectors satisfying $a_i + b_i = c_i + d_i$ for each $i = 1, \dots, m$. Then*

$$d_{\text{hel}}^2(M_a, M_b) = d_{\text{hel}}^2(M_c, M_d).$$

Proof We claim the following analogue of Lemma 11.3.18: for any $X_1^m = x_1^m$ and any communication transcript τ , we may write

$$Q(\Pi(x_1^m) = \tau \mid x_1^m) = \prod_{i=1}^m f_{i, x_i}(\tau) \quad (11.5.3)$$

for some functions f_{i, x_i} . Indeed, letting $\tau = \{z_i^{(t)}\}_{i \leq n, t \leq T}$ we have

$$Q(\Pi(x_1^m) = \tau \mid x_1^m) = \prod_{i,t} Q(z_i^{(t)} \mid x_1^m, z_{\rightarrow i}^{(t)}) = \prod_{i=1}^m \underbrace{\prod_{t=1}^T Q(z_i^{(t)} \mid x_i, z_{\rightarrow i}^{(t)})}_{=: f_{i, x_i}(\tau)}$$

where we use that message $z_i^{(t)}$ depends only on x_i and $z_{\rightarrow i}^{(t)}$. Then we can write $M_b(\Pi(X_1^m) = \tau)$ as a product using Eq. (11.5.3): integrating over independent $X_i \sim P_{b_i}$, we have

$$M_b(\Pi(X_1^m) = \tau) = \int Q(\tau \mid x_1^m) dP_{b_1}(x_1) \cdots dP_{b_m}(x_m) = \prod_{i=1}^m \underbrace{\int f_{i, \tau}(x_i) dP_{b_i}(x_i)}_{=: g_{i, b_i}(\tau)} = \prod_{i=1}^m g_{i, b_i}(\tau).$$

Taking M_a, M_b, M_c, M_d as in the statement of the lemma,

$$d_{\text{hel}}^2(M_a, M_b) = 1 - \sum_{\tau} \sqrt{\prod_{i=1}^m g_{i,a_i}(\tau) g_{i,b_i}(\tau)}.$$

But as $a_i + b_i = c_i + d_i$ and each is $\{0, 1\}$ -valued, we certainly have $g_{i,a_i} g_{i,b_i} = g_{i,c_i} g_{i,d_i}$, and so the lemma follows. \square

The second result we require is due to Jayram [120], and is the following:

Lemma 11.5.3. *Let $\{P_b\}_{b \in \{0,1\}^m}$ be any collection of distributions satisfying the cutting and pasting property $d_{\text{hel}}^2(P_a, P_b) = d_{\text{hel}}^2(P_c, P_d)$ whenever $a, b, c, d \in \{0,1\}^m$ satisfy $a + b = c + d$. Let $N = 2^k$ for some $k \in \mathbb{N}$. Then for any collection of bit vectors $\{b^{(i)}\}_{i=1}^N \subset \{0,1\}^m$ with $\langle b^{(i)}, b^{(j)} \rangle = 0$ for all $i \neq j$ and $b = \sum_i b^{(i)}$,*

$$\prod_{l=1}^k (1 - 2^{-l}) d_{\text{hel}}^2(P_0, P_b) \leq \sum_{i=1}^m d_{\text{hel}}^2(P_0, P_{b^{(i)}}).$$

We defer the technical proof to Section 11.5.1.

A computation shows that $\prod_{l=1}^k (1 - 2^{-l}) > \frac{2}{7}$. Lemma 11.5.3 nearly gives us our desired result (11.5.2), except that it requires a power of 2. To that end, let k_0 be the largest $k \in \mathbb{N}$ such that $2^{k_0} \leq m$, and construct bit vectors $b^{(1)}, \dots, b^{(2^{k_0})}$ satisfying $\sum_i b^{(i)} = \mathbf{1}$ and $1 \leq \|b^{(i)}\|_0 \leq 2$ for each i . Then Lemma 11.5.3, via the cutting-pasting property of the marginals M , implies

$$\frac{2}{7} d_{\text{hel}}^2(M_0, M_1) \leq \sum_{i=1}^{2^{k_0}} d_{\text{hel}}^2(M_0, M_{b^{(i)}}) \leq 2 \sum_{i=1}^m d_{\text{hel}}^2(M_0, M_{e_i}),$$

where the second inequality again follows from Lemma 11.5.3 as $b^{(i)} = e_j$ or $e_j + e_{j'}$ for some basis vectors $e_j, e_{j'}$. This gives Lemma 11.5.1. \square

Step 2: from Hellinger to Shannon information

Now we relate the strong data processing constants for mutual information in Definition 11.7 to compare Hellinger distances with mutual information. We claim the following lemma.

Lemma 11.5.4. *Let the conditions of Theorem 11.4.4 hold. Let M_0 and M_{e_l} be the marginal distributions over Π when X_i have the sampling distribution (11.5.1). Then for $l \in \{1, \dots, m\}$,*

$$d_{\text{hel}}^2(M_{e_l}, M_0) \leq \frac{c+1}{2} \beta I(X_l; \Pi(X_1^m) \mid V = 0).$$

Proof Consider the following alternative distributions. Let $W \sim \text{Uniform}\{0, 1\}$, and draw $X' \in \mathcal{X}^m$ with independent coordinates according to

$$X'_i \stackrel{\text{iid}}{\sim} P_0 \text{ if } W = 0 \quad \text{or} \quad X'_i \sim \begin{cases} P_0 & \text{if } i \neq l \\ P_1 & \text{if } i = l \end{cases} \text{ if } W = 1.$$

Then we have the Markov chain $W \rightarrow X' \rightarrow \Pi(X')$, and moreover,

$$W \rightarrow X'_l \rightarrow \Pi(X') \leftarrow X'_l,$$

so that additionally $W \rightarrow X'_l \rightarrow \Pi(X')$ is a Markov chain. As a consequence, Definition 11.7 of the strong data processing inequality gives

$$I(W; \Pi(X')) \leq \beta I(X'_l; \Pi(X')).$$

Using Proposition 2.2.10, we thus have

$$d_{\text{hel}}^2(M_{e_l}, M_0) \leq I(W; \Pi(X')) \leq \beta I(X'_l; \Pi(X')). \quad (11.5.4)$$

It remains to relate $I(X'_l; \Pi(X'))$ to $I(X_l; \Pi(X) \mid V = 0)$. Here we bounded likelihood ratio between P_0 by P_1 . Indeed, we have by the condition (11.4.3) that

$$P_0 \geq \frac{1}{c} P_1 \quad \text{so} \quad (c+1)P_0 \geq P_0 + P_1 \quad \text{or} \quad P_0 \geq \frac{2}{c+1} \frac{P_0 + P_1}{2}.$$

As a consequence, we have

$$\begin{aligned} I(X_l; \Pi(X'_l) \mid V = 0) &= \int D_{\text{kl}}(Q(\cdot \mid X_l = x) \parallel M_0) dP_0(x) \\ &\geq \frac{2}{c+1} \int D_{\text{kl}}(Q(\cdot \mid X_l = x) \parallel M_0) \frac{dP_0(x) + dP_1(x)}{2} \\ &\geq \frac{2}{c+1} \int D_{\text{kl}}(Q(\cdot \mid X_l = x) \parallel \bar{M}) \frac{dP_0(x) + dP_1(x)}{2} \\ &= \frac{2}{c+1} I(X'_l; \Pi(X')), \end{aligned}$$

where the second inequality uses that $\bar{M} = \int Q(\cdot \mid X_l = x) \frac{dP_0(x) + dP_1(x)}{2}$ minimizes the integrated KL-divergence (recall inequality (10.2.3)). Returning to inequality (11.5.4), we evidently have the result of the lemma. \square

Step 3: Completing the proof of Theorem 11.4.4

By combining the tensorization Lemma 11.5.1 with the information bound in Lemma 11.5.4, we obtain

$$d_{\text{hel}}^2(M_0, M_1) \leq 7 \sum_{i=1}^m d_{\text{hel}}^2(M_0, M_{e_i}) \leq \frac{7}{2}(c+1)\beta \sum_{i=1}^m I(X_i; \Pi \mid V = 0).$$

By symmetry, we also have

$$d_{\text{hel}}^2(M_0, M_1) \leq 7 \sum_{i=1}^m d_{\text{hel}}^2(M_0, M_{e_i}) \leq \frac{7}{2}(c+1)\beta \sum_{i=1}^m I(X_i; \Pi \mid V = 1).$$

Now, we note that as the X_i are independent conditional on V (and w.l.o.g. for the purposes of mutual information, we may assume they are discrete), for any $v \in \{0, 1\}$ we have

$$\begin{aligned}
 \sum_{i=1}^m I(X_i; \Pi \mid V = v) &= \sum_{i=1}^m [H(X_i \mid V = v) - H(X_i \mid \Pi, V = v)] \\
 &= \sum_{i=1}^m [H(X_i \mid X_1^{i-1}, V = v) - H(X_i \mid \Pi, V = v)] \\
 &\leq \sum_{i=1}^m [H(X_i \mid X_1^{i-1}, V = v) - H(X_i \mid X_1^{i-1}, \Pi, V = v)] \\
 &= \sum_{i=1}^m I(X_i; \Pi \mid X_1^{i-1}, V = v) = I(X_1, \dots, X_m; \Pi \mid V = v),
 \end{aligned}$$

where the inequality used that conditioning decreases entropy. We thus obtain

$$d_{\text{hel}}^2(M_0, M_1) \leq \frac{7}{2}(c+1)\beta \min_{v \in \{0,1\}} I(X_1, \dots, X_m; \Pi \mid V = v)$$

as desired.

11.5.1 Proof of Lemma 11.5.3

We prove the result by induction. It is trivially true for $m = 1$, that is, $k = 0$, so now we consider the inductive case, that is, it holds for $m = 1, \dots, 2^{k-1}$ and we consider $m = 2^k$.

First, we observe that if $\{u_i\}_{i=1}^N$ are arbitrary vectors and $\bar{u} = \frac{1}{N} \sum_{i=1}^N u_i$ is their mean, then

$$\sum_{i,j} \|u_i - u_j\|_2^2 = \sum_{i,j} \|u_i - \bar{u} + \bar{u} - u_j\|_2^2 = \sum_{i,j} \|u_i - \bar{u}\|_2^2 + \sum_{i,j} \|\bar{u} - u_j\|_2^2 = 2N \sum_{i=1}^N \|\bar{u} - u_i\|_2^2.$$

Thus, if u_0 is any other vector, that \bar{u} minimizes $\sum_i \|u_i - u\|_2^2$ over all u gives

$$\frac{1}{N} \sum_{1 \leq i < j \leq N} \|u_i - u_j\|_2^2 \leq \sum_{i=1}^N \|u_i - \bar{u}\|_2^2 \leq \sum_{i=1}^N \|u_i - u_0\|_2^2. \quad (11.5.5)$$

Now, we return to the Hellinger distances. Evidently $2d_{\text{hel}}^2(P_a, P_b) = \|\sqrt{p_a}(\cdot) - \sqrt{p_b}(\cdot)\|_2^2$, so that it is a Euclidean distance. As a consequence, for any pairwise disjoint collection of N bit vectors $b^{(i)}$, we have

$$\sum_{i=1}^N d_{\text{hel}}^2(P_0, P_{b^{(i)}}) \geq \frac{1}{N} \sum_{1 \leq i < j \leq N} d_{\text{hel}}^2(P_{b^{(i)}}, P_{b^{(j)}}) = \frac{1}{N} \sum_{1 \leq i < j \leq N} d_{\text{hel}}^2(P_0, P_{b^{(i)} + b^{(j)}}) \quad (11.5.6)$$

where the inequality follows from (11.5.5) and the equality by the assumed cut-and-paste property. Now, we apply Baranyai's theorem, which says that we may decompose any complete graph K_N , where N is even, into $N - 1$ perfect matchings \mathcal{M}_i with $N/2$ edges—necessarily, as they form a

perfect matching—where each \mathcal{M}_i is edge disjoint. Identifying the pairs $i < j$ with the complete graph, we thus obtain

$$\sum_{1 \leq i < j \leq N} d_{\text{hel}}^2(P_0, P_{b^{(i)}+b^{(j)}}) = \sum_{l=1}^{N-1} \sum_{(i,j) \in \mathcal{M}_l} d_{\text{hel}}^2(P_0, P_{b^{(i)}+b^{(j)}}). \quad (11.5.7)$$

Now fix $n \in \{1, \dots, N-1\}$ and a matching \mathcal{M}_n . By assumption we have $\langle b^{(i)}+b^{(j)}, b^{(i')}+b^{(j')} \rangle = 0$ for any distinct pairs $(i, j), (i', j') \in \mathcal{M}_n$, and moreover, $\sum_{(i,j) \in \mathcal{M}_n} (b^{(i)} + b^{(j)}) = b$. Thus, our induction hypothesis gives that for any $l \in \{1, \dots, N-1\}$ and any of our matchings \mathcal{M}_n , we have

$$\sum_{(i,j) \in \mathcal{M}_n} d_{\text{hel}}^2(P_0, P_{b^{(i)}+b^{(j)}}) \geq d_{\text{hel}}^2(P_0, P_b) \prod_{l=1}^{k-1} (1 - 2^{-l}).$$

Substituting this lower bound into inequality (11.5.7) and using inequality (11.5.6), we obtain

$$\sum_{i=1}^N d_{\text{hel}}^2(P_0, P_{b^{(i)}}) \geq \frac{1}{N} \cdot (N-1) d_{\text{hel}}^2(P_0, P_b) \prod_{l=1}^{k-1} (1 - 2^{-l}) = d_{\text{hel}}^2(P_0, P_b) \prod_{l=1}^k (1 - 2^{-l}),$$

where we have used $N = 2^k$.

11.6 Bibliography

Data processing inequalities originate with Dobrushin’s study of central limit theorems for Markov chains [66, 67]; Dobrushin first proved Proposition 11.1.1 (see [67, Sec. 3.1]). Cohen et al. [54] show that the strong data processing constant for variation distance is the largest of the strong data processing constants (Theorem 11.1.2) for finite state spaces using careful linear algebraic techniques, also showing the opposite extremality (inequality (11.1.1)) of the χ^2 contraction coefficient [54, Proposition II.6.15] for finite state spaces. Del Moral et al. [65] and Polyanskiy and Wu [154] give related and approachable treatments for general alphabets, and Exercises 11.1 and 11.2 follow [65]. More broadly, strong data processing inequalities arise in many applications in communication, estimation, and some functional analysis [157, 154].

Communication complexity begins with Yao [193], which introduces the communication complexity setting we discuss in Section 11.3, making the connections between randomized complexities and public (shared) randomness. The standard classical reference for the subject is Kushilevitz and Nisan’s book [129]. There are numerous techniques that we do not discuss, including so-called discrepancy lower bounds, which address both randomized and deterministic communication complexity; for example, these give the stronger lower bound that $\text{DCC}_\delta(\text{IP}_2) \geq n - O(1)$ [129, Example 3.29 and Exercise 3.30]. Communication complexity has uses far beyond the “standard” communication setting we have outlined, with more recent research showing how to use the techniques to provide lower bounds on the performance of algorithms in many computational models, such as streaming models and memory-limited computation [149, 159]. Our information complexity approach follows Bar-Yossef et al. [15]. Recent work has shown how communication lower bounds and strong data processing inequalities can be used to show the necessity of “memorization” in some natural problems in machine learning, where any learning procedure with good enough performance necessarily encodes substantial irrelevant information about a dataset [40].

Our treatment of communication complexity and its applications in estimation follows an approach Zhang et al. [198] originate. The particular techniques we adapt, involving direct sums and strong data processing in communication, we adapt from Braverman et al. [39] and Garg et al. [99]. Our results apply most easily to scenarios in which each machine or agent owns only a single data item, which allows application of Proposition 11.4.2; tensorizing this to multiple observations requires some care, but can be done with a truncation argument [198, 39] or more careful Sobolev inequalities [157]. Our extension to private estimation scenarios follows the paper [70], which also shows how to generalize to other variants of privacy.

11.7 Exercises

Exercise 11.1 (Approximating nonnegative convex functions): Let $f : \mathbb{R} \rightarrow \mathbb{R}_+ \cup \{+\infty\}$ be a closed, nonnegative convex function.

- Show that there exists a sequence of piecewise linear functions f_n satisfying $f_{n-1} \leq f_n \leq f$ for all n and for which $f_n(x) \uparrow f(x)$ pointwise for all x s.t. $f(x) < \infty$, and $f_n(x) \uparrow \infty$ otherwise. *Hint:* Let \mathcal{L} be the collection of linear functions below f , that is $\mathcal{L} = \{l \mid l(x) = a + bx, l(x) \leq f(x) \text{ for all } x\}$, and note that $f(x) = \sup\{l(x) \mid l \in \mathcal{L}\}$. (See Appendix C.2.) You may replace \mathcal{L} with functions of the form $l(x) = f(x_0) + g(x - x_0)$, where $g \in \partial f(x_0)$ is a subderivative of f at x_0 .
- Show that if for some $z_0 \in \mathbb{R}$ we have $f(z_0) = 0$, then one may take the functions f_n to be of the form $f_n(x) = \sum_{i=1}^n a_i [b_i - x]_+ + \sum_{i=1}^n c_i [x - d_i]_+$, where $b_i \leq z_0$, $d_i \geq z_0$, and $a_i, c_i \geq 0$.
- Conclude that for any measure μ on \mathbb{R}_+ , $\int f_n d\mu \uparrow \int f d\mu$.

Exercise 11.2 (Proving Theorem 11.1.2): In this question, we formalize the sketched proof of Theorem 11.1.2 by filling in details of the following steps. Let $\alpha = \alpha_{\text{TV}}(Q)$ be the Dobrushin coefficient of the channel Q and $f : \mathbb{R} \rightarrow \mathbb{R}_+ \cup \{+\infty\}$ be a closed convex function.

- There exists a nondecreasing sequence f_n of piecewise linear functions, each of the form $f_n(x) = \sum_{i=1}^n a_i [b_i - x]_+ + \sum_{i=1}^n c_i [x - d_i]_+$, where $b_i \leq 1$, $d_i \geq 1$, and $a_i, c_i \geq 0$. *Hint:* Exercise 11.1.
- Let $M_v(A) = \int Q(A \mid x) dP_v(x)$ for $v \in \{0, 1\}$ be the induced marginal distributions. Show that for any function of the form $h(t) = [t - \Delta]_+$, where $\Delta > 1$,

$$D_h(M_0 \| M_1) \leq \alpha D_h(P_0 \| P_1) \quad (11.7.1)$$

by the following steps:

- Define the set $\mathcal{X}(\Delta) := \{x \mid dP_0(x) \leq \Delta dP_1(x)\}$. Argue that $\mathcal{X}(\Delta)$ must be non-null (i.e., have positive measure).
- Define the probability distribution P_Δ with density

$$dP_\Delta(x) = \frac{\Delta dP_1(x) - dP_0(x)}{\int [\Delta dP_1(x) - dP_0(x)]_+} \mathbf{1}_{\{x \in \mathcal{X}(\Delta)\}}.$$

Argue that the measure

$$G = \Delta P_1 - (\Delta - 1)P_\Delta$$

is a probability distribution.

iii. Show that

$$D_h(P_0 \| P_1) = \|P_0 - G\|_{\text{TV}}.$$

It may be useful to show that $dP_0 - dG \leq 0$ on $\mathcal{X}(\Delta)$.

iv. Conclude that

$$D_h(P_0 \| P_1) \geq \frac{1}{\alpha} \|Q \circ P_0 - Q \circ G\|_{\text{TV}} \geq \frac{1}{\alpha} D_h(Q \circ P_0 \| Q \circ P_1).$$

(c) Using the monotone convergence theorem, show that $D_f(M_0 \| M_1) \leq \alpha D_f(P_0 \| P_1)$.

Exercise 11.3 (Markov chain mixing): Consider a Markov chain X_1, X_2, \dots with transition distribution $P(\cdot | x)$ and stationary distribution π . Let $P^k(\cdot | x)$ denote the distribution of the Markov chain initialized in state x after k steps. Assume there exists some (finite) positive integer $k \in \mathbb{N}$ such that for any two initial states x_0, x_1 , the Markov chain satisfies

$$\left\| P^k(\cdot | x_0) - P^k(\cdot | x_1) \right\|_{\text{TV}} \leq \beta < 1.$$

Show that the Markov chain enjoys *fast mixing* for any f divergence: if there is any n such that $D_f(P^n(\cdot | x) \| \pi) < \infty$, the Markov chain mixes exponentially quickly in that it satisfies

$$\limsup_n \frac{1}{n} \log D_f(P^n(\cdot | x) \| \pi) \leq \frac{1}{k} \log \beta < 0.$$

In brief, as soon as one can demonstrate a constant gap in variation distance, one is guaranteed a Markov chain mixes geometrically.

Exercise 11.4: For $k \in [1, \infty]$, we consider the collection of distributions

$$\mathcal{P}_k := \{P : \mathbb{E}_P[|X|^k]^{1/k} \leq 1\},$$

that is, distributions P supported on \mathbb{R} with k th moment bounded by 1. We consider minimax estimation of the mean $\mathbb{E}[X]$ for these families under ε -local differential privacy, meaning that for each observation X_i , we observe a private realization Z_i (which may depend on Z_1^{i-1}) where Z_i is an ε -differentially private view of X_i . Let \mathcal{Q}_ε denote the collection of all ε -differentially private channels, and define the (locally) private minimax risk

$$\mathfrak{M}_n(\theta(\mathcal{P}), (\cdot)^2, \varepsilon) := \inf_{\hat{\theta}_n} \inf_{Q \in \mathcal{Q}_\varepsilon} \sup_{P \in \mathcal{P}} \mathbb{E}_{P, Q}[(\hat{\theta}_n(Z_1^n) - \theta(P))^2].$$

(a) Assume that $\varepsilon \leq 1$. For $k \in [1, \infty]$, show that there exists a constant $c > 0$ such that

$$\mathfrak{M}_n(\theta(\mathcal{P}_k), (\cdot)^2, \varepsilon) \geq c \left(\frac{1}{n\varepsilon^2} \right)^{\frac{k-1}{k}}.$$

(b) Give an ε -locally differentially private estimator achieving the minimax rate in part (a).

Exercise 11.5: Show that strong data processing inequality in Theorem 11.2.1 is sharp in the following sense. There exist ε -differentially private channels Q_ε such that for any Bernoulli distributions P_0 and P_1 and induced marginal distributions $M_{v, \varepsilon} = Q(\cdot | X = 1)P_v(X = 1) + Q(\cdot | X = 0)P_v(X = 0)$,

$$\frac{D_{\text{kl}}(M_{0, \varepsilon} \| M_{1, \varepsilon})}{\|P_0 - P_1\|_{\text{TV}}^2} = \frac{\varepsilon^2}{2} + O(\varepsilon^3)$$

as $\varepsilon \downarrow 0$.

Exercise 11.6: We apply the results of Exercise 11.4 to a problem of estimation of drug use. Assume we interview a series of individuals $i = 1, \dots, n$, asking whether each takes illicit drugs. Let $X_i \in \{0, 1\}$ be 1 if person i uses drugs, 0 otherwise, and define $\theta^* = \mathbb{E}[X] = \mathbb{E}[X_i] = P(X = 1)$. Instead of X_i we observe answers Z_i under differential privacy,

$$Z_i \mid X_i = x \sim Q(\cdot \mid X_i = x)$$

for a ε -differentially private Q with $\varepsilon < \frac{1}{2}$ (so that $(e^\varepsilon - 1)^2 \leq 2\varepsilon^2$). Let \mathcal{Q}_ε denote the family of all ε -differentially private channels, and let \mathcal{P} denote the Bernoulli distributions with parameter $\theta(P) = P(X_i = 1) \in [0, 1]$ for $P \in \mathcal{P}$.

- (a) Use Le Cam's method and the strong data processing inequality in Theorem 11.2.1 to show that the minimax rate for estimation of the proportion θ^* in absolute value satisfies

$$\mathfrak{M}_n(\theta(\mathcal{P}), |\cdot|, \varepsilon) := \inf_{Q \in \mathcal{Q}_\varepsilon} \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_{P, Q} \left[|\hat{\theta}(Z_1, \dots, Z_n) - \theta(P)| \right] \geq c \frac{1}{\sqrt{n\varepsilon^2}},$$

where $c > 0$ is a universal constant.

- (b) Give a rate-optimal estimator for this problem. That is, define an ε -differentially private channel Q and an estimator $\hat{\theta}$ such that $\mathbb{E}[|\hat{\theta}(Z_1^n) - \theta|] \leq C/\sqrt{n\varepsilon^2}$, where C is a universal constant.
- (c) Download the dataset at <http://web.stanford.edu/class/stats311/Data/drugs.txt>, which consists of a sample of 100,000 hospital admissions and whether the patient was abusing drugs (a 1 indicates abuse, 0 no abuse). Use your estimator from part (b) to estimate the population proportion of drug abusers: give an estimated number of users for $\varepsilon \in \{2^{-k}, k = 0, 1, \dots, 10\}$. Perform each experiment several times. Assuming that the proportion of users in the dataset is the true population proportion, how accurate is your estimator?

Exercise 11.7: Show that the randomized communication complexity (11.3.1) satisfies $\text{RCC}_\delta(f) \leq O(1) \log_{\frac{1}{\delta}} \text{RCC}_{1/3}(f)$ for any f and any $\delta < 1$.

Exercise 11.8 (From public to private randomness): Consider the randomized complexity (11.3.1) and associated public-randomness complexity $\text{RCC}_\delta^{\text{pub}}$. Let $\mathcal{X} = \mathcal{Y} = \{0, 1\}^n$ and $f : \mathcal{X} \times \mathcal{Y} \rightarrow \{0, 1\}$, and let Π be a protocol using public randomness U such that $\max_{x, y} \mathbb{P}(\Pi(x, y, U) \neq f(x, y)) \leq \epsilon$.

- (a) Use Hoeffding's inequality to show that there are $k = \frac{\log 2}{\delta^2} n$ points u_1, \dots, u_k such that if $I \in [k]$ is chosen uniformly at random, then $\mathbb{P}(\Pi(x, y, u_I) \neq f(x, y)) \leq \epsilon + \delta$.
- (b) Give a protocol that uses no public randomness but whose communication complexity is at most $\text{depth}(\Pi) + O(1) \log \frac{n}{\delta}$.
- (c) Conclude that $\text{RCC}_\delta(f) \leq \text{RCC}_\delta^{\text{pub}}(f) + O(1) \log \frac{n}{\delta}$.

Exercise 11.9 (An information lower bound for indexing): In the indexing problem in communication complexity, Alice receives an n -bit string $x \in \{0, 1\}^n$ and Bob an index $y \in [n] = \{1, \dots, n\}$, and the two communicate to evaluate x_y ; set $f(x, y) = x_y$.

- (a) Show that if Bob can send messages, the communication complexity of indexing satisfies $\text{CC}(f) \leq O(1) \log n$.

In the *one way* communication model, only Alice can send messages. Let μ be the uniform distribution on $(X, Y) \in \{0, 1\}^n \times [n]$. We will show that $\text{DCC}_\delta^\mu(f) \geq (1 - h_2(\delta))n$, where $h_2(p) = -p \log_2 p - (1 - p) \log_2 (1 - p)$ is the binary entropy.

- (b) Fix the index $Y = i$ and let $p_i = \mathbb{P}(\hat{X}_i = X_i \mid Y = i)$ based on a protocol Π . Use Fano's inequality (Proposition 9.4.1) to argue that $h_2(p_i) \geq H_2(X_i \mid \Pi)$.
- (c) Show that if Π is a δ -error one-way protocol under μ , then

$$I(X_1^n; \Pi) \geq (1 - h_2(\delta))n.$$

Exercise 11.10 (Information complexity for entrywise less or equal): Consider the entrywise less than or equal to function $f : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}$ with $f(x, y) = \mathbf{1}\{x \preceq y\}$, so that $f(x, y) = 1$ if $x_i \leq y_i$ for each i and 0 if there exists i such that $x_i > y_i$.

- (a) Show that f has the decompositional structure (11.3.4). Give the functions g and h .
- (b) Give a fooling distribution μ on $\mathcal{X} \times \mathcal{Y}$ for f .
- (c) Use Theorem 11.3.13 and a modification of the proof of Proposition 11.3.15 to show that $\text{IC}_\delta(f) \geq \frac{n}{4}(1 - 2\sqrt{\delta(1 - \delta)})$. (This is order optimal, because $\text{IC}_\delta(f) \leq \text{CC}(f) \leq n + 1$ trivially.)

Exercise 11.11 (Lower bounds for private logistic regression): This question is (likely) challenging. Consider the logistic regression model for $y \in \{\pm 1\}$, $x \in \mathbb{R}^d$, that

$$p_\theta(y \mid x) = \frac{1}{1 + \exp(-y\langle \theta, x \rangle)}.$$

For a distribution P on $(X, Y) \in \mathbb{R}^d \times \{\pm 1\}$, where $Y \mid X = x$ has logistic distribution, define the excess risk

$$r(\theta, P) := \mathbb{E}_P[\ell(\theta; X, Y)] - \inf_{\theta} \mathbb{E}_P[\ell(\theta; X, Y)]$$

where $\ell(\theta; x, y) = \log(1 + \exp(-y\langle x, \theta \rangle))$ is the logistic loss. Let \mathcal{P} be the collection of such distributions, where X is supported on $\{-1, 1\}^d$. Peeking ahead to Chapter 17 and Section 17.3, for a channel Q mapping $(X, Y) \rightarrow Z$, define

$$\mathfrak{M}_n(\mathcal{P}, Q) := \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_{P, Q}[r(\hat{\theta}(Z_1^n), P)],$$

where the expectation is taken over $Z_i \sim Q(\cdot \mid X_i, Z_1^{i-1})$. Assume that the channel releases are all (locally) ε -differentially private.

- (a) Show that for all n large enough,

$$\mathfrak{M}_n(\mathcal{P}, Q) \geq c \cdot \frac{d}{n} \cdot \frac{d}{\varepsilon \wedge \varepsilon^2}$$

for some (numerical) constant $c > 0$.

- (b) Suppose we allow additional passes through the dataset (i.e. multiple rounds of communication), but still require that all data Z_i released from X_i be ε -differentially private. That is, assume we have the (sequential and interactive) release schemes of Fig. 11.3, and we guarantee that

$$Z_i^{(t)} \sim Q(\cdot \mid X_i, B^{(1)}, \dots, B^{(t)}, Z_1^{(t)}, \dots, Z_{i-1}^{(t)})$$

is $\varepsilon_{i,t}$ -differentially private, where $\sum_t \varepsilon_{i,t} \leq \varepsilon$ for all i . Does the lower bound of part (a) change?

Exercise 11.12: In this question, we prove Proposition 11.4.2.

- (a) Show that if $p(v)$ and $p(v \mid x)$ denote the p.m.f.s of V and V conditional on $X = x$, then

$$e^{-\alpha} p(v) \leq p(v \mid x) \leq e^{\alpha} p(v).$$

- (b) Show that

$$|p(v \mid z) - p(v)| \leq 2(e^{\alpha} - 1) \|P_X(\cdot \mid z) - P_X(\cdot)\|_{\text{TV}}.$$

- (c) Complete the proof of the proposition.

JCD Comment: A few additional exercises to add:

1. Prove Yao's minimax theorem.
2. Is there a clean "memorization" phenomenon to cover?

Chapter 12

Squared error and asymptotically exact optimality guarantees

JCD Comment: Add a bit on exact optimality I guess

JCD Comment: Notation for derivative matrices?

The squared error takes a particularly central place in the theory of estimation, as its particularly simple structure means it admits closed forms for many optimal estimators, and, for lack of a better description, it plays nicely with differentiation. In this chapter, we develop a few of the elements of this theory, presenting some of the classical bounds on estimation accuracy, extending them beyond basic parametric estimators. For example, in any Bayesian problem, where we draw a parameter θ from a prior π on θ and then observe a sample $X \sim P_\theta$, the posterior mean $\mathbb{E}[\theta \mid X]$ is always Bayes optimal: we have

$$\mathbb{E} \left[\|\hat{\theta}(X) - \theta\|_2^2 \right] = \mathbb{E} \left[\mathbb{E} \left[\|\hat{\theta}(X) - \theta\|_2^2 \mid X \right] \right] \geq \mathbb{E} \left[\mathbb{E} \left[\|\mathbb{E}[\theta \mid X] - \theta\|_2^2 \right] \right],$$

where we use that $\inf_t \mathbb{E}[\|Y - t\|_2^2] = \mathbb{E}[\|Y - \mathbb{E}[Y]\|_2^2]$ for any random vector Y .

Rather explicitly using this Bayesian perspective with its explicit characterization of optimal estimation, however, we leverage the connection between squared error and correlation, showing that on average over parameters θ , the parameter θ is correlated with the *score* vector

$$\dot{\ell}_\theta(x) := \nabla_\theta \log p_\theta(x) = \frac{\nabla_\theta p_\theta(x)}{p_\theta(x)}$$

when $X \sim P_\theta$ has density p_θ . The connection between correlation—as an inner product—and squared error then allows us to provide strong lower bounds on estimation. Another description for our approaches here is that integration-by-parts implies lower bounds on the squared error.

The main results in this chapter provide lower bounds on the estimation error in terms of the *Fisher information matrix* of the parameter θ ,

$$J(\theta) := \mathbb{E}_\theta \left[\dot{\ell}_\theta(X) \dot{\ell}_\theta(X)^\top \right] = \text{Cov}_\theta(\dot{\ell}_\theta(X)), \quad (12.0.1)$$

where the equality exploits that, so long as we may exchange integration and differentiation,

$$\mathbb{E}_\theta[\dot{\ell}_\theta(X)] = \int p_\theta(x) \frac{\nabla p_\theta(x)}{p_\theta(x)} dx = \int \nabla p_\theta(x) dx \stackrel{(*)}{=} \nabla_\theta \int p_\theta(x) dx = 0,$$

because $\int p_\theta = 1$ for any θ , where equality (\star) typically requires justification. For an i.i.d. sample $X_1^n \stackrel{\text{iid}}{\sim} P_\theta$, then, the n -observation information matrix satisfies

$$J_n(\theta) = nJ(\theta)$$

because $\log p_\theta(x_1^n) = \sum_{i=1}^n \log p_\theta(x_i)$. Then the main consequences of the development here are inequalities of the form

$$\mathbb{E} \left[\|\hat{\theta}(X_1^n) - \theta\|_2^2 \right] \geq \frac{\text{tr}(J(\theta)^{-1})}{n} - O(1/n^2),$$

so that the Fisher information matrix (12.0.1) plays a fundamental role in estimation lower bounds.

12.1 The Cramér-Rao inequality

To maintain some connection with historical approaches, we begin with the Cramér-Rao bound, which provides lower bounds on the mean-squared error of *unbiased* estimators. We will be a bit fast-and-loose here when consider the integrability conditions, as we will not use the bound to develop any actual optimality results. Nonetheless, we provide the (heuristic) setting here: for a parameter $\theta \in \mathbb{R}$, assume that X has density p_θ with respect to some base measure μ , where $\dot{p}_\theta(x) := \frac{\partial}{\partial \theta} p_\theta(x)$ exists for all x and is appropriately integrable.

Proposition 12.1.1 (The Cramér-Rao bound). *Assume that $\hat{\theta} : \mathcal{X} \rightarrow \mathbb{R}$ is unbiased for θ . Then*

$$\mathbb{E}_\theta \left[(\hat{\theta}(X) - \theta)^2 \right] \geq \frac{1}{J(\theta)}$$

under appropriate conditions on the density p_θ .

Proof We use integration by parts, that is, that $\int u dv = uv - \int v du$. Because $\mathbb{E}_\theta[\hat{\theta}] = \theta$, we have

$$1 = \frac{\partial}{\partial \theta} \mathbb{E}_\theta[\hat{\theta}] = \frac{\partial}{\partial \theta} \int \hat{\theta}(x) p_\theta(x) d\mu(x) \stackrel{(\star)}{=} \int \hat{\theta}(x) \dot{p}_\theta(x) d\mu(x),$$

where (\star) exchanges integration and differentiation. Then as $\dot{p}_\theta(x) = \dot{\ell}_\theta(x) p_\theta(x)$, we obtain

$$1 = \int \hat{\theta}(x) \dot{\ell}_\theta(x) p_\theta(x) d\mu(x) = \mathbb{E}_\theta[\hat{\theta}(X) \dot{\ell}_\theta(X)] = \mathbb{E}_\theta \left[(\hat{\theta}(X) - \theta) \dot{\ell}_\theta(X) \right]$$

because $\mathbb{E}_\theta[\dot{\ell}_\theta(X)] = 0$. Applying Cauchy-Schwarz to observe $\mathbb{E}_\theta[(\hat{\theta}(X) - \theta) \dot{\ell}_\theta(X)]^2 \leq \mathbb{E}_\theta[(\hat{\theta}(X) - \theta)^2] \mathbb{E}_\theta[\dot{\ell}_\theta^2(X)]$ gives the result. \square

While there is a long history of using the Cramér-Rao bound to claim some type of statistical efficiency for an estimator, the Cramér-Rao bound is the biggest con in the history of statistics. In no way does it actually indicate that an estimator is efficient, as most practically efficient estimators are biased, because of regularization or other choices. Even if an estimator is asymptotically unbiased, e.g., $\mathbb{E}_n[\hat{\theta}_n(X_1^n)] \rightarrow \theta$, the Cramér-Rao inequality says nothing. Additionally, soon as the underlying parameter space Θ is compact, typically, no unbiased estimator that guarantees $\hat{\theta} \in \Theta$ even *exists*.

12.1.1 Compact sets and the failure of the Cramér-Rao bound

We justify the preceding polemical screed via exemplar results. For the first, consider estimation of bounded parameters. We begin with estimation on an interval.

Example 12.1.2 (Estimation of a bounded parameter): Let $\{P_\theta\}$ be any family of probability models absolutely continuous with respect to one another, meaning that for any $\theta_0, \theta_1 \in \Theta$, $P_{\theta_0}(A) > 0$ implies $P_{\theta_1}(A) > 0$. Suppose the parameter $\theta \in [a, b]$, and assume $\hat{\theta}$ is unbiased for θ but still satisfies $\hat{\theta} \in [a, b]$. Then we claim that $\hat{\theta} = a$ with probability 1 and $\hat{\theta} = b$ with probability 1. Indeed, when $\theta = a$, we have

$$\mathbb{E}_a[\hat{\theta}(X)] = a,$$

and as $\hat{\theta} \geq a$, we must have $\hat{\theta}(X) = a$ with probability 1. Similarly, for $\theta = b$, we have $\hat{\theta} \leq b$ and as $\mathbb{E}_b[\hat{\theta}(X)] = b$ then $\hat{\theta}(X) = b$ with probability 1. This is obviously a contradiction unless $a = b$. \diamond

Building off of example 12.1.2, we can show that if we have the *a priori* guarantee that the parameter $\theta \in \Theta$ for a convex body Θ , then we can reduce the mean-squared error of any estimator $\hat{\theta}$ trivially: simply project $\hat{\theta}$ onto Θ . We can provide stricter guarantees. For example, Exercise 12.1 shows that if $\hat{\theta}$ is unbiased for θ , then (except in pathological settings) the estimator

$$\tilde{\theta}(X) := \text{Proj}_\Theta(\hat{\theta}) = \underset{\theta \in \Theta}{\operatorname{argmin}} \|\theta - \hat{\theta}\|_2^2$$

has strictly lower mean-squared error

$$\mathbb{E}_\theta \left[\|\tilde{\theta}(X) - \theta\|_2^2 \right] < \mathbb{E}_\theta \left[\|\hat{\theta}(X) - \theta\|_2^2 \right]$$

for *all* $\theta \in \Theta$: a rebuke to the idea that the Cramér-Rao inequality provides a fundamental limit.

12.1.2 Regularization and the failure of the Cramér-Rao bound

An alternative perspective arises when we consider regularized estimators, which are related to projection estimators but can be analytically simpler. Ridge or ℓ_2 -regularization provides the most common example; for example, for linear regression, this is

$$\hat{\theta}_\lambda = \underset{\theta}{\operatorname{argmin}} \left\{ \|X\theta - Y\|_2^2 + \lambda \|\theta\|_2^2 \right\}.$$

In general, there always exists $\lambda > 0$ with smaller mean-squared error than the unregularized estimator $\hat{\theta}_0$ (see Exercise 12.2).

JCD Comment: Give the Fisher information in this problem?

Here, let us consider the somewhat estimation of the mean of a standard normal random vector. Let $Z \sim \mathcal{N}(\theta, I_d)$, where $\theta \in \mathbb{R}^d$ is unknown, and note that the Fisher information in this case is $J(\theta) = d$. Note that

$$\underset{t}{\operatorname{argmin}} \left\{ \|t - Z\|_2^2 + \lambda \|t\|_2^2 \right\} = \frac{Z}{1 + \lambda}$$

takes the simple shrinkage form $Z \mapsto Z/(1 + \lambda)$, we reparameterize and consider estimators of the form

$$\hat{\theta}_\beta = \beta Z$$

for $\beta \in [0, 1]$. (This is equivalent to the choice $\lambda = \frac{1}{\beta} - 1$ in the usual ridge regression formulation). Then immediately we see that as $Z = \theta + \varepsilon$ for $\varepsilon \sim \mathbf{N}(0, I_d)$, we have $\hat{\theta}_\beta = \beta\theta + \beta\varepsilon$, and so

$$\mathbb{E} \left[\|\hat{\theta}_\beta - \theta\|_2^2 \right] = (1 - \beta)^2 \|\theta\|_2^2 + \beta^2 d.$$

Taking derivatives, we see by inspection that $\beta = \frac{\|\theta\|_2^2}{d + \|\theta\|_2^2}$ minimizes the mean-squared error. This at least shows that there is always *some* shrinkage-based estimator outperforming the putative information bound, which is d .

In the particular case of normal mean estimation, James-Stein-type estimators allow us to say more. As we wish to estimate θ , the “optimal” shrinkage parameter $\beta = \frac{\|\theta\|_2^2}{d + \|\theta\|_2^2}$ is of course unavailable. But $\mathbb{E}[\|Z\|_2^2] = d + \|\theta\|_2^2$, and so an estimate of β is possible: we have $\beta = \frac{\mathbb{E}[\|Z\|_2^2] - d}{\mathbb{E}[\|Z\|_2^2]} = 1 - \frac{d}{\mathbb{E}[\|Z\|_2^2]}$ and replacing β with a the slightly less conservative counterpart

$$\hat{\beta} := 1 - \frac{[d - 2]_+}{\|Z\|_2^2}$$

gives the *James-Stein* estimator

$$\hat{\theta}_{\text{JS}} := \left(1 - \frac{[d - 2]_+}{\|Z\|_2^2} \right) Z.$$

Famously, we have the following result.

Theorem 12.1.3. *Let the dimension $d > 2$ and $Z \sim \mathbf{N}(\theta, I_d)$ for some $\theta \in \mathbb{R}^d$. Then*

$$\mathbb{E}_\theta \left[\|\hat{\theta}_{\text{JS}} - \theta\|_2^2 \right] < \mathbb{E}_\theta \left[\|Z - \theta\|_2^2 \right] = d.$$

In this case, we see that an adaptive estimator strictly outperforms the sample mean, which achieves the information bound.

JCD Comment: Give citation and commentary for James-Stein, saying that we won’t worry about it too much, just that it means we ought to move beyond Cramér-Rao.

12.2 The van Trees inequality: a Bayesian Cramér-Rao bound

While our discussion in the preceding sections casts a dim view on the Fisher information “efficiency” bound that the Cramér-Rao inequality implies, it turns out that its main failing was in attempting to provide a pointwise bound valid for any parameter θ . By a bit of Bayesian averaging—as we do in essentially all of our minimax lower bounds—we can address many of the issues and show that similar information bounds hold, at least as the sample size n grows. The key will be that if we put a prior π on the parameter θ of interest, then *on average* over π , the parameter θ must be correlated with the score $\dot{\ell}_\theta = \nabla \log p_\theta$.

12.2.1 The van Trees inequality in one dimension

We first develop the so-called *van Trees* inequality in one-dimensional problems. Let $\theta \in \mathbb{R}$ be the parameter of interest, and assume that $\{p_\theta\}$ is a family of densities with respect to some base measure μ . Assume that the prior density π has support $[a, b]$, where $\pi(a) = \pi(b) = 0$, and that it is differentiable with finite prior Fisher information

$$J(\pi) := \int_a^b \frac{\pi'(t)^2}{\pi(t)} dt = \int_a^b ((\log \pi(t))')^2 \pi(t) dt.$$

Then in this case, we can show that a Fisher information lower bound holds *on average* over π , with a small additional penalty to capture the information π gives on the parameter θ . For notational simplicity in the theorem, we assume that $X \sim P_\theta$, detailing the i.i.d. case afterward.

Theorem 12.2.1 (The van Trees inequality). *Let the above conditions hold. Then*

$$\mathbb{E}_\pi \left[\mathbb{E}_\theta \left[(\hat{\theta}(X) - \theta)^2 \right] \right] \geq \frac{1}{\mathbb{E}_\pi[J(\theta)] + J(\pi)}.$$

Proof Define the augmented score

$$\dot{\ell}_{\theta,\pi}(x) := \frac{\partial}{\partial \theta} \log(p_\theta(x)\pi(\theta)) = \frac{(p_\theta(x)\pi(\theta))'}{p_\theta(x)\pi(\theta)} = \frac{\dot{p}_\theta(x)}{p_\theta(x)} + \frac{\dot{\pi}(\theta)}{\pi(\theta)}.$$

The key is that $\hat{\theta}(X) - \theta$ and $\dot{\ell}_{\theta,\pi}(X)$ are correlated, which then implies a lower bound. Taking an expectation jointly over $\theta \sim \pi$ and $X \sim p_\theta$, we can apply integration-by-parts (with $u = \hat{\theta}(x) - \theta$ and $dv = (\pi(\theta)p_\theta(x))' d\theta$ in the identity $\int u dv = uv - \int v du$) to obtain

$$\begin{aligned} \mathbb{E} \left[(\hat{\theta}(X) - \theta) \dot{\ell}_{\theta,\pi}(X) \right] &= \iint (\hat{\theta}(x) - \theta) \dot{\ell}_{\theta,\pi}(x) \pi(\theta) p_\theta(x) d\mu(x) d\theta \\ &= \iint (\hat{\theta}(x) - \theta) (\pi(\theta) p_\theta(x))' d\theta d\mu(x) \\ &= \int \left(\left[(\hat{\theta}(x) - \theta) \pi(\theta) p_\theta(x) \right]_a^b + \int \pi(\theta) p_\theta(x) d\theta \right) d\mu(x) \\ &= 1, \end{aligned}$$

because $\pi(a) = \pi(b) = 0$. So applying the Cauchy-Schwarz inequality, we have

$$\mathbb{E} \left[(\hat{\theta}(X) - \theta)^2 \right]^{1/2} \mathbb{E} \left[\dot{\ell}_{\theta,\pi}(X)^2 \right]^{1/2} \geq 1.$$

Finally, we note that

$$\mathbb{E}[\dot{\ell}_{\theta,\pi}(X)^2] = \mathbb{E}[(\dot{\ell}_\theta(X) + \pi'(\theta)/\pi(\theta))^2] = \mathbb{E}_\pi[J(\theta)] + J(\pi),$$

because $\mathbb{E}_\pi[\pi'(\theta)/\pi(\theta)] = 0$. □

When the sample X is an i.i.d. sample $X_1^n \stackrel{\text{iid}}{\sim} P_\theta$, then of course the Fisher information becomes $J_n(\theta) = nJ(\theta) = n\mathbb{E}_\theta[\dot{\ell}_\theta(X)^2]$, where X is a single observation. Then Theorem 12.2.1 provides the bound

$$\mathbb{E} \left[(\hat{\theta}(X_1^n) - \theta)^2 \right] \geq \frac{1}{n\mathbb{E}_\pi[J(\theta)] + J(\pi)} = \frac{1}{n\mathbb{E}_\pi[J(\theta)]} - O(1) \frac{J(\pi)}{\mathbb{E}_\pi[J(\theta)]^2 n^2},$$

assuming $J(\pi) < \infty$. In principle it is possible to choose the prior π to maximize the lower bound, but typically this is rather immaterial.

Because we usually want a *local* lower bound—something akin to the Cramér-Rao bound that (almost) applies to a particular problem parameter θ_0 of interest—it is most common to localize the prior π to some shrinking neighborhood of θ_0 . Let π be any prior supported on a neighborhood of 0; for simplicity, we can take $[-1, 1]$. Then for any $a > 0$ and $\theta_0 \in \mathbb{R}$, the density

$$f_a(\theta) := \frac{1}{a} \cdot \pi\left(\frac{\theta - \theta_0}{a}\right)$$

satisfies $\int f_a(\theta) d\theta = 1$, has support $[-a, a]$, and information $J(f_a) = \frac{1}{a} J(\pi)$. (See also Exercise 12.3.) Then if $J(\theta)$ is continuous in θ near θ_0 , we obtain the following corollary.

Corollary 12.2.2. *Let $\theta_0 \in \mathbb{R}$, the prior density π have compact support on $[-1, 1]$ and finite information $J(\pi) < \infty$, and $a_n > 0$ be any sequence satisfying $\frac{1}{n} \ll a_n \ll 1$. Define the prior densities*

$$\pi_n(\theta) := \frac{1}{a_n} \cdot \pi\left(\frac{\theta - \theta_0}{a_n}\right).$$

If $J(\theta)$ is continuous in a neighborhood of θ_0 , then

$$\int_{\theta_0 - a_n}^{\theta_0 + a_n} \mathbb{E}_\theta \left[(\hat{\theta}(X_1^n) - \theta)^2 \right] \pi_n(\theta) d\theta \geq \frac{1}{nJ(\theta_0)} - o(1/n).$$

In a strong sense, then, the inverse Fisher information $J(\theta)^{-1}$ provides a (nearly) pointwise lower bound on estimation in mean-square error.

12.2.2 The van Trees inequality in d -dimensions

We can extend Theorem 12.2.1 to higher-dimensional scenarios as well: more or less, we simply apply the preceding lower bound to each coordinate of the parameter. In this case, instead of a scalar Fisher information, we have information matrices

$$J(\theta) := \mathbb{E}_\theta[\dot{\ell}_\theta(X)\dot{\ell}_\theta(X)^\top] \quad \text{and} \quad J(\pi) = \mathbb{E}_\pi[\nabla \log \pi(\theta) \nabla \log \pi(\theta)^\top].$$

With this notation, the following multivariate van Trees inequality follows almost exactly the same lines of proof as Theorem 12.2.1.

Theorem 12.2.3 (Basic multivariate van Trees). *Let $\pi = \pi_1 \times \cdots \times \pi_d$ be a product distribution on $\theta \in \mathbb{R}^d$, where each coordinate density π_j has support $[a_j, b_j]$ with $\pi_j(a_j) = \pi_j(b_j)$. Then for any matrix $A \succ 0$,*

$$\mathbb{E}_\pi \left[\mathbb{E}_\theta [(\hat{\theta}(X) - \theta)^\top A^{-1} (\hat{\theta}(X) - \theta)] \right] \geq \frac{d^2}{\mathbb{E}_\pi [\text{tr}(AJ(\theta))] + \text{tr}(AJ(\pi))}.$$

Proof Let $\pi(\theta) = \prod_{j=1}^d \pi_j(\theta_j)$, and let $\pi_{\setminus j}(\theta_{\setminus j}) = \prod_{k \neq j} \pi_k(\theta_k)$. For shorthand, let $\partial_j = \frac{\partial}{\partial \theta_j}$.

Expanding in each coordinate, we have

$$\begin{aligned}
\langle \widehat{\theta}(x) - \theta, \dot{\ell}_{\theta, \pi}(x) \rangle &= \sum_{j=1}^d (\widehat{\theta}_j(x) - \theta_j) \frac{\partial}{\partial \theta_j} \log(p_\theta(x) \pi(\theta)) \\
&= \sum_{j=1}^d (\widehat{\theta}_j(x) - \theta_j) \left([\dot{\ell}_\theta(x)]_j + \frac{\pi'_j(\theta_j)}{\pi_j(\theta_j)} \right) \\
&= \sum_{j=1}^d (\widehat{\theta}_j(x) - \theta_j) \frac{\partial_j(p_\theta(x) \pi_j(\theta_j))}{p_\theta(x) \pi_j(\theta_j)}.
\end{aligned}$$

Then each term satisfies integration-by-parts identities as in the proof of Theorem 12.2.1, with

$$\begin{aligned}
\mathbb{E} \left[(\widehat{\theta}_j(X) - \theta_j) \cdot [\dot{\ell}_{\theta, \pi}(X)]_j \right] &= \iint (\widehat{\theta}_j(x) - \theta_j) \partial_j(p_\theta(x) \pi_j(x)) d\theta_j \cdot \pi_{\setminus j}(\theta_{\setminus j}) d\theta_{\setminus j} d\mu(x) \\
&= \iint \left([(\widehat{\theta}_j(x) - \theta_j) p_\theta(x) \pi_j(\theta_j)]_{\theta_j=a_j}^{\theta_j=b_j} + \int_{a_j}^{b_j} \pi_j(\theta_j) p_\theta(x) d\theta_j \right) \pi_{\setminus j}(\theta_{\setminus j}) d\theta_{\setminus j} d\mu(x) \\
&= \iint \pi(\theta) p_\theta(x) d\theta d\mu(x) = 1.
\end{aligned}$$

We therefore obtain

$$\mathbb{E} \left[\langle \widehat{\theta}(X) - \theta, \dot{\ell}_{\theta, \pi}(X) \rangle \right] = \iint \langle \widehat{\theta}(x) - \theta, \dot{\ell}_{\theta, \pi}(x) \rangle \pi(\theta) p_\theta(x) d\mu(x) d\theta = d.$$

Applying the Cauchy-Schwarz inequality to the inner product $\langle u, v \rangle = \langle A^{-1/2}u, A^{1/2}v \rangle$ for any $A \succ 0$ gives

$$\langle \widehat{\theta} - \theta, \dot{\ell}_{\theta, \pi} \rangle = \langle \widehat{\theta} - \theta, \dot{\ell}_\theta \rangle + \langle \widehat{\theta} - \theta, \nabla \log \pi(\theta) \rangle \leq \|\widehat{\theta} - \theta\|_{A^{-1}} \left(\left\| \dot{\ell}_\theta + \nabla \log \pi(\theta) \right\|_A \right),$$

and applying Cauchy-Schwarz and recognizing that

$$\mathbb{E}_\pi \left[\left\| \dot{\ell}_\theta(X) + \nabla \log \pi(\theta) \right\|_A^2 \right] = \mathbb{E}_\pi [\text{tr}(AJ(\theta))] + \text{tr}(AJ(\pi))$$

gives the theorem. \square

Typically, we choose the matrix A to be related to the Fisher information of the parameter θ . As motivation, recall our (heuristic) development in Chapter 3, where we showed in the approximation (3.2.4) that

$$\widehat{\theta}_n - \theta^* \sim \mathbf{N}(0, n^{-1} J(\theta^*)^{-1}),$$

where we have substituted the Fisher information. Recalling Chapter 5.3.4, a small modification of Corollary 5.3.11 and its application of Corollary 5.3.10 yields the following result. (See Exercise 12.4.)

Corollary 12.2.4. *Let the conditions of Corollary 5.3.11 hold, let π be any density supported on a convex body Θ_0 , and let Θ be any convex body with $\Theta \supset \Theta_0$. Let $\widehat{\theta}_n$ be the maximum likelihood estimator and define the projected estimator*

$$\widetilde{\theta}_n = \text{Proj}_\Theta(\widehat{\theta}_n).$$

Then for any matrix $B(\theta) \succeq 0$ continuous in θ , uniformly in $\theta \in \Theta_0$

$$\mathbb{E}_\theta \left[\|\tilde{\theta}_n(X_1^n) - \theta\|_{B(\theta)}^2 \right] \leq \frac{1}{n} \text{tr}(B(\theta) J(\theta)^{-1}) + o(1/n).$$

The most natural idea to link Theorem 12.2.3 with Corollary 12.2.4 is to make the estimator $\hat{\theta}_n$ “pivotal,” meaning that its asymptotic distribution is independent of θ^* . Then (again using our heuristics), we multiply by $J(\theta^*)^{1/2}$, so that for $X_1^n \stackrel{\text{iid}}{\sim} P_{\theta^*}$,

$$J(\theta^*)^{1/2} (\hat{\theta}_n - \theta^*) \sim \mathbf{N}(0, n^{-1} I_d).$$

Let us assume $J(\theta)$ is continuous in θ . In Corollary 12.2.4, setting $B(\theta) = J(\theta)$ we then obtain

$$\mathbb{E}_\theta \left[(\tilde{\theta}_n(X_1^n) - \theta)^\top J(\theta) (\tilde{\theta}_n(X_1^n) - \theta) \right] \leq \frac{d}{n} + o(1/n)$$

simultaneously for all $\theta \in \Theta_0$. Comparing with Theorem 12.2.3, fix θ^* . Let $\delta > 0$ and Θ_0 to be any neighborhood of θ^* small enough that $(1 - \delta)J(\theta) \preceq J(\theta^*) \preceq (1 + \delta)J(\theta)$ for all $\theta \in \Theta_0$. Then

$$\mathbb{E}_\pi \left[\|\hat{\theta}_n(X_1^n) - \theta\|_{A^{-1}}^2 \right] \geq \frac{d^2}{n \text{tr}(A \mathbb{E}_\pi[J(\theta)]) + \text{tr}(AJ(\pi))} \geq \frac{d^2}{n(1 + \delta) \text{tr}(AJ(\theta^*)) + \text{tr}(AJ(\pi))}$$

for any estimator $\hat{\theta}_n$. Take $A = J(\theta^*)^{-1}$ to obtain

$$\mathbb{E}_\pi \left[(\hat{\theta}_n(X_1^n) - \theta)^\top J(\theta^*) (\hat{\theta}_n(X_1^n) - \theta) \right] \geq \frac{d^2}{n(1 + \delta)d + \text{tr}(J(\theta^*)^{-1}J(\pi))} = \frac{d}{(1 + \delta)n} - O(1/n^2),$$

showing that the van Trees inequality is essentially unimprovable.

12.2.3 The van Trees inequality for a function of the parameter

We now present two results on the van Trees inequality in higher dimensions. In many statistical problems, it is interesting to estimate a function of a parameter of the underlying distribution rather than a parameter identifying the distribution itself. So suppose we have a statistic $T : \mathbb{R}^d \rightarrow \mathbb{R}^p$ and wish to estimate the parameter $T(\theta)$ rather than θ itself. The next example highlights a simple case where $T(\theta) = e_j^\top \theta$ for a standard basis vector e_j .

Example 12.2.5 (Treatment effect estimation): Consider a randomized controlled trial, where individuals are either assigned to treatment (with an indicator variable $Z = 1$) or control (with indicator variable $Z = 0$). We use the *potential outcomes* framework to let the response Y of an individual be $Y(0)$ under no treatment and $Y(1)$ under treatment, recognizing that for any individual we observe only one of $Y(0)$ and $Y(1)$. Then we model the response linearly by

$$Y = Y(Z) = \beta_0 + tZ + \beta^\top X + \varepsilon,$$

where β_0 is an intercept and $X \in \mathbb{R}^d$ are covariates, so the $\theta = (\beta_0, t, \beta) \in \mathbb{R}^{d+2}$ is the vector of parameters. Assuming that $Z \sim \text{Uniform}\{0, 1\}$, independently of X and $(Y(0), Y(1))$, $\theta^* = \text{argmin}_\theta \mathbb{E}[(Y - \beta_0 - tZ - \beta^\top X)^2]$ has t -component $t^* = \mathbb{E}[Y(1)] - \mathbb{E}[Y(0)]$. This value is the *average treatment effect*. (See Exercise 12.5.) \diamond

In Section 12.3, we exploit these ideas profitably in general estimation problems, such as M-estimation without well-specified models or the general estimation lower bound framework we develop in Chapter 9.

In any case, suppose we have a statistic $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^p$ and wish to estimate $\psi(\theta)$. Assume the derivative matrix $\dot{\psi}(\theta) \in \mathbb{R}^{p \times d}$ exists, with entries $\dot{\psi}_{ij}(\theta) = \frac{\partial}{\partial \theta_j} \psi_i(\theta)$, so that $\psi(\theta + \Delta) = \psi(\theta) + \dot{\psi}(\theta)\Delta + o(\|\Delta\|)$. We then have the following theorem, which extends Theorem 12.2.3.

Theorem 12.2.6. *Let the conditions of Theorem 12.2.3 hold, and let $C \in \mathbb{R}^{d \times p}$ be an arbitrary matrix. Assume that ψ is continuously differentiable on the support of π . Then for any $A \succ 0$ and any estimator $\hat{\psi}$ of $\psi(\theta)$,*

$$\mathbb{E}_\pi \left[(\hat{\psi}(X) - \psi(\theta))^\top A^{-1} (\hat{\psi}(X) - \psi(\theta)) \right] \geq \frac{\mathbb{E}_\pi [\text{tr}(\dot{\psi}(\theta)C)^2]}{\text{tr}(CAC^\top \mathbb{E}_\pi[J(\theta)]) + \text{tr}(CAC^\top J(\pi))}.$$

Proof Let $C \in \mathbb{R}^{d \times p}$ have rows c_j^\top , so $C = [c_1 \ \cdots \ c_d]^\top$, and define the error vector $E(x, \theta) = \hat{\psi}(x) - \psi(\theta)$. Recognize that

$$C^\top \dot{\ell}_{\theta, \pi}(x) = \sum_{j=1}^d c_j \partial_j \log(p_\theta(x) \pi(\theta)) = \sum_{j=1}^d c_j \left(\frac{\partial_j p_\theta(x)}{p_\theta(x)} + \frac{\partial_j \pi_j(\theta_j)}{\pi_j(\theta_j)} \right).$$

Then we have

$$\begin{aligned} \mathbb{E}[\langle \hat{\psi}(X) - \psi(\theta), C^\top \dot{\ell}_{\theta, \pi}(X) \rangle] &= \sum_{j=1}^d \iint \langle E(x, \theta), c_j \rangle \partial_j (p_\theta(x) \pi_j(\theta_j)) d\theta_j \pi_{\setminus j}(\theta_{\setminus j}) d\theta_{\setminus j} d\mu(x) \\ &= \sum_{j=1}^d \iint \left([\langle E(x, \theta), c_j \rangle p_\theta(x) \pi_j(\theta_j)]_{a_j}^{b_j} + \langle \partial_j \psi(\theta), c_j \rangle p_\theta(x) \pi_j(\theta_j) \right) \pi_{\setminus j}(\theta_{\setminus j}) d\mu(x) \\ &= \mathbb{E}_\pi \left[\text{tr}(\dot{\psi}(\theta)C) \right] \end{aligned}$$

by integration by parts.

For the remainder of the proof, we mimic that of Theorem 12.2.3. We have

$$\mathbb{E} \left[\langle E(X, \theta), C^\top \dot{\ell}_{\theta, \pi}(X) \rangle \right] \leq \mathbb{E} \left[\|E(X, \theta)\|_{A^{-1}}^2 \right]^{1/2} \mathbb{E} \left[\|C^\top \dot{\ell}_{\theta, \pi}(X)\|_A^2 \right]^{1/2},$$

and because the scores are mean zero, we also have

$$\mathbb{E}_\pi \left[\|C^\top \dot{\ell}_{\theta, \pi}(X)\|^2 \right] = \mathbb{E}_\pi \left[\|C^\top \dot{\ell}_\theta(X) + \nabla \log \pi(\theta)\|_A^2 \right] = \mathbb{E}_\pi \left[\text{tr}(CAC^\top J(\theta)) \right] + \text{tr}(CAC^\top J(\pi)).$$

Squaring and dividing through gives the result. \square

The final version of the van Trees, or Bayesian Cramér-Rao inequality, extends Theorem 12.2.3 to include additional terms to better reflect the geometry of the problem at hand, which eliminates the need for some of the local approximations we have made to consider particular parameters θ_0

of interest. In this case, we let the matrices $A \succ 0$ and $C \in \mathbb{R}^{d \times p}$ vary in θ in Theorem 12.2.6. For such a setting, we will require a prior Fisher information with

$$s_{C,\pi}(\theta) := \frac{1}{\pi(\theta)} \left[\sum_{i=1}^d \frac{\partial}{\partial \theta_i} (C_{ij}(\theta) \pi(\theta)) \right]_{j=1}^p \in \mathbb{R}^p$$

taking the place of the usual score vector $\nabla \log \pi(\theta)$ (these agree when $p = d$ and $C = I_d$, so that $s_{I,\pi}(\theta) = \nabla \log \pi(\theta)$) and

$$J(\pi \mid A, C) := \mathbb{E}_\pi \left[s_{C,\pi}(\theta)^\top A(\theta) s_{C,\pi}(\theta) \right]. \quad (12.2.1)$$

To state the most general extension of Theorem 12.2.6, we will present a few additional definitions and conditions. A function f on $\Theta \subset \mathbb{R}^d$ is *suitably continuous* if for each coordinate j and almost all values $\theta \in \Theta$, the coordinate function $h(t) := f(\theta + te_j)$ is absolutely continuous. Then we consider the conditions (i)–(v) below, which relax the assumptions necessary for Theorem 12.2.6.

- (i) The density $p_\theta(x)$ is measurable in (x, θ) and, for almost every x , is suitably continuous in θ .
- (ii) $\psi : \Theta \rightarrow \mathbb{R}^p$ and $C : \Theta \rightarrow \mathbb{R}^{d \times p}$ are suitably continuous.
- (iii) The Fisher information $J(\theta) := \mathbb{E}_\theta[\dot{\ell}_\theta(X) \dot{\ell}_\theta(X)^\top]$ exists and $\text{diag}(J(\theta))^{1/2}$ is locally integrable.
- (iv) $A(\theta)$ is positive definite and continuous in θ .
- (v) The prior density π is suitably continuous and its domain Θ is compact with piecewise \mathcal{C}^1 boundary, and $\pi(\text{int } \Theta) > 0$ and $\pi(\text{bd } \Theta) = \{0\}$.

Then Gill and Levit [102, Theorem 1] show the following result.

Theorem 12.2.7 (The multivariate van Trees inequality). *Let conditions (i)–(v) hold. Then for any estimator $\hat{\psi} : \mathcal{X} \rightarrow \mathbb{R}^p$,*

$$\mathbb{E}_\pi \left[(\hat{\psi}(X) - \psi(\theta))^\top A(\theta)^{-1} (\hat{\psi}(X) - \psi(\theta)) \right] \geq \frac{\mathbb{E}_\pi[\text{tr}(\dot{\psi}(\theta) C(\theta))]^2}{\mathbb{E}_\pi[\text{tr}(C(\theta) A(\theta) C(\theta)^\top J(\theta))] + J(\pi \mid A, C)}.$$

Exercise 12.6 asks you to prove a slightly weaker version of Theorem 12.2.7, which follows from arguments similar to those we use to prove Theorem 12.2.6. Under the same conditions, when we have an i.i.d. sample $X_1^n \stackrel{\text{iid}}{\sim} P_\theta$, that the Fisher information tensorizes as well implies

Corollary 12.2.8. *Let the conditions of Theorem 12.2.7 hold and $X_1^n \stackrel{\text{iid}}{\sim} P_\theta$. Then for any estimator $\hat{\psi}_n : \mathcal{X}^n \rightarrow \mathbb{R}^p$,*

$$\mathbb{E}_\pi \left[(\hat{\psi}_n(X_1^n) - \psi(\theta))^\top A(\theta)^{-1} (\hat{\psi}_n(X_1^n) - \psi(\theta)) \right] \geq \frac{\mathbb{E}_\pi[\text{tr}(\dot{\psi}(\theta) C(\theta))]^2}{n \mathbb{E}_\pi[\text{tr}(C(\theta) A(\theta) C(\theta)^\top J(\theta))] + J(\pi \mid A, C)}.$$

We provide a corollaries to Theorem 12.2.7 by choosing the matrices A and C appropriately; typically, we take them to be related to the Fisher information matrix $J(\theta)$ as in our discussion after Corollary 12.2.4. We begin with the “natural” parameterization, where we take ψ to be the identity and C the identity, and let $A(\theta)^{-1} = J(\theta)$ be the Fisher information. Then so long as $J(\theta) \succ 0$ and is continuous, using $\frac{1}{n+c} \geq \frac{1}{n} - \frac{c}{n^2}$, we have the following corollary.

Corollary 12.2.9. *Let the conditions of Theorem 12.2.7 hold. Then for $A(\theta) = J(\theta)$ and $C(\theta) = I_d$, we have $J(\pi | A, C) = \mathbb{E}_\pi[\nabla \log \pi(\theta) \nabla \log \pi(\theta)^\top]$ and*

$$\mathbb{E}_\pi \left[(\hat{\theta}_n(X_1^n) - \theta)^\top J(\theta) (\hat{\theta}_n(X_1^n) - \theta) \right] \geq \frac{d^2}{nd + J(\pi | A, C)} \geq \frac{d}{n} - \frac{J(\pi | A, C)}{n^2}.$$

Taking ψ to be the identity once again, but this time letting $A(\theta) = I_d$ and $C(\theta) = J(\theta)^{-1}$ and assuming J is suitably differentiable, we obtain

$$\mathbb{E}_\pi \left[\|\hat{\theta}_n(X_1^n) - \theta\|_2^2 \right] \geq \frac{\mathbb{E}_\pi[\text{tr}(J(\theta)^{-1})]^2}{n \mathbb{E}_\pi[\text{tr}(J(\theta)^{-1})] + J(\pi | C)} \geq \frac{\mathbb{E}_\pi[\text{tr}(J(\theta)^{-1})]}{n} - \frac{J(\pi | C)}{n^2}.$$

When $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^p$ is not the identity mapping, the bounds are more sophisticated, but natural choices again present themselves. We begin by heuristically developing an upper bound, using a technique known as the *delta method*: as ψ is assumed differentiable, for $\hat{\theta}_n$ obeying the usual asymptotics $\hat{\theta}_n - \theta^* \sim \mathbf{N}(0, (1/n)\Sigma)$ for some covariance Σ , we may proceed heuristically to obtain

$$\begin{aligned} \psi(\hat{\theta}_n) - \psi(\theta^*) &= \dot{\psi}(\theta^*)(\hat{\theta}_n - \theta^*) + \underbrace{O(\|\hat{\theta}_n - \theta^*\|^2)}_{=O(1/n)} \\ &\dot{\sim} \mathbf{N} \left(0, \frac{1}{n} \dot{\psi}(\theta^*) \Sigma \dot{\psi}(\theta^*)^\top \right). \end{aligned}$$

(Any text on asymptotic statistics provides rigorous versions of these claims, for example, van der Vaart [185, Ch. 3]. See also Exercise 5.12.)

Assume that $p \leq d$ and that $\dot{\psi}(\theta) \in \mathbb{R}^{p \times d}$ is rank p , so that the Fisher information for $\psi(\theta)$ in the model P_θ is, by a change of variables,

$$J_\psi(\theta) := \left(\dot{\psi}(\theta) J(\theta)^{-1} \dot{\psi}(\theta)^\top \right)^{-1} \in \mathbb{R}^{p \times p}. \quad (12.2.2)$$

To obtain an analogue of the dimension-dependent bounds in Corollary 12.2.9, we can therefore take A to be the inverse information for ψ , $A(\theta) = J_\psi(\theta)^{-1}$, and $C(\theta) = J(\theta)^{-1} \dot{\psi}(\theta)^\top J_\psi(\theta) \in \mathbb{R}^{d \times p}$. Then

$$\dot{\psi}(\theta) C(\theta) = \dot{\psi}(\theta) J(\theta)^{-1} \dot{\psi}(\theta)^\top J_\psi(\theta) = I_p,$$

while the cyclic property of the trace and that $C(\theta) A(\theta) = J(\theta)^{-1} \dot{\psi}(\theta)^\top$ imply

$$\text{tr}(C(\theta) A(\theta) C(\theta)^\top J_n(\theta)) = n \text{tr} \left(J(\theta)^{-1} \dot{\psi}(\theta)^\top J_\psi(\theta) \dot{\psi}(\theta) \right) = n \text{tr}(I_p) = np.$$

Letting $\tilde{J}(\pi | \psi) := J(\pi | A, C)$ as in (12.2.1) for these choices of A and C , we then obtain the following “normalized” risk inequality.

Corollary 12.2.10. *In addition to the conditions of Theorem 12.2.7, assume that $\dot{\psi}(\theta) J(\theta)^{-1} \dot{\psi}(\theta)^\top$ is continuously differentiable in θ and positive definite. Then letting $J_\psi(\theta)$ be the Fisher information (12.2.2) for ψ in the model $\{P_\theta\}$ satisfies*

$$\mathbb{E}_\pi \left[(\hat{\psi}_n(X_1^n) - \psi(\theta))^\top J_\psi(\theta) (\hat{\psi}_n(X_1^n) - \psi(\theta)) \right] \geq \frac{p^2}{np + \tilde{J}(\pi | \psi)} \geq \frac{p}{n} - \frac{\tilde{J}(\pi | \psi)}{n^2}.$$

Let us revisit linear regression (or Example 12.2.5) in this context, where we assume that the model is true.

Example 12.2.11: Assume the linear regression model

$$Y_i = X_i^\top \theta + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathbf{N}(0, \sigma^2).$$

Because the log-likelihood is $\log p_\theta(y \mid x) = -\frac{1}{2\sigma^2}(x_i^\top \theta - y_i)^2 - \frac{1}{2} \log(2\pi\sigma^2)$, the Fisher information for the parameter vector θ given n observations (X_i, Y_i) in this model is

$$J_n(\theta) = \sum_{i=1}^n \frac{1}{\sigma^4} \mathbb{E}[\varepsilon_i X_i X_i^\top \varepsilon_i] = \frac{1}{\sigma^2} \sum_{i=1}^n X_i X_i^\top = \frac{1}{\sigma^2} X^\top X,$$

where we let $X = [X_1 \cdots X_n]^\top$ be the design matrix. If we wish to estimate a single coordinate $\psi(\theta) = e_j^\top \theta$, the information for the coordinate θ_j is then $J_\psi(\theta) = \frac{1}{\sigma^2} (e_j^\top (X^\top X)^{-1} e_j)^{-1}$, and Corollary 12.2.10 gives the lower bound

$$\mathbb{E}_\pi [(\hat{\theta}_j - \theta_j)^2] \geq \sigma^2 e_j^\top (X^\top X)^{-1} e_j + O(1/n^2) = \frac{\sigma^2}{n} e_j^\top \left(n^{-1} X^\top X \right)^{-1} e_j + O(1/n^2).$$

Because the standard regression estimator $\hat{\theta} = (X^\top X)^{-1} X^\top Y = \theta + (X^\top X)^{-1} X^\top \varepsilon$, we have $\mathbb{E}[(\hat{\theta}_j - \theta_j)^2] = \sigma^2 e_j^\top (X^\top X)^{-1} e_j$, achieving the lower bound. \diamond

The lower bounds in Corollary 12.2.10 provide further insights in cases where we must estimate problems with nuisance parameters that are uninteresting, or at least not the subject of particular scientific investigation. See Exercise 12.7.

12.3 Beyond parametric problems

We often have statistical problems in which there is no particular model underlying the problem, and we wish only to estimate a parameter defined in terms of the underlying distribution generating the data. For example, in the M-estimation problems in Chapter 5.3, the quantity of interest is

$$\theta(P) := \operatorname{argmin}_{\theta} \mathbb{E}_P[\ell(\theta, Z)],$$

which is a function of the distribution P . Treating the probability distribution P as the “parameter,” then $\theta(P)$ is a function of that parameter, which at the outset seems hard to address with the techniques we have developed thus far, as in this case P is (typically) infinite dimensional. This is also the more abstract setting we adopt in Chapter 9.

To extend the techniques for lower bounding estimator error to this setting, we adopt a perspective that begins with Stein [172], where estimation in a general problem should be at least as hard as estimation in any particular (model-based) sub-problem. Thus, we take the following two-phase approach to developing estimation lower bounds for a parameter $\theta(P)$ around a base distribution P_0 :

1. For the base distribution, define a model family $\{P_t\}$ indexed by $t \in \mathbb{R}^k$
2. Show how to write $\theta(P_t)$ as a function $\theta(t)$ of the underlying parameter t , then apply Theorem 12.2.7 or its corollaries.
3. Optionally, revisit step 1 and modify the underlying models to make estimation error the largest

Exercises 9.10, 9.12, and 9.13 explore this approach in the context of information-theoretic lower bounds. Here, we employ the ideas to obtain asymptotically exact results.

The actual approach to doing this is, at the end of the day, rather straightforward given our development of minimax lower bounds as well as the information-based bounds in the preceding sections. Assume without loss of generality that P_0 has a density p_0 on \mathcal{X} with respect to a base measure μ (we could always simply take $\mu = P_0$). Then for a bounded function $g : \mathcal{X} \rightarrow \mathbb{R}^k$ with mean $\mathbb{E}_0[g(X)] = 0$, define the *tilted density*

$$p_t(x) := (1 + \langle t, g(x) \rangle) p_0(x). \quad (12.3.1)$$

It is immediate that for any $t \in \mathbb{R}^k$, we have $\int p_t(x) d\mu(x) = \mathbb{E}_0[(1 + \langle t, g(X) \rangle)] = 1$, and because we assume g is bounded, for t near enough 0 we are guaranteed that $p_t \geq 0$, so that p_t is indeed a probability density. Once we have these tilted densities, we will assume that the parameter $\theta(P)$ of interest is locally differentiable:

Assumption A.12.1. *Let U be any neighborhood of 0. For the model family $\{P_t\}_{t \in U}$, the parameter $\theta(P_t)$ is continuously differentiable in t on U , with derivative matrix $\dot{\theta}(t) \in \mathbb{R}^{d \times k}$.*

Assumption A.12.1 is abstract, but frequently holds; we provide an example with mean estimation here, and after presenting the main lower bound leveraging the assumption, show how it applies to the general M-estimation problems of Chapter 5.3.

Example 12.3.1 (Nonparametric mean estimation): For distributions P on $X \in \mathbb{R}^d$, define $\theta(P) := \mathbb{E}_P[X]$. Then for the tilted family $\{P_t\}$ of (12.3.1),

$$\theta(P_t) = \mathbb{E}_0[(1 + \langle t, g(X) \rangle)X] = \mathbb{E}_0[X] + \text{Cov}_0(X, g(X))t,$$

where we have used that g is mean zero. Then $\theta(P_t)$ is affine in t and hence differentiable. \diamond

Conveniently, the tilted construction (12.3.1) immediately allows us to compute the Fisher information matrix for the nuisance parameter t we have invented. Because $\log p_t(x) = \log(1 + \langle t, g(x) \rangle) + \log p_0(x)$, we obtain

$$\nabla \log p_t(x) = \frac{g(x)}{1 + \langle t, g(x) \rangle} \quad \text{and} \quad \mathbb{E}_t[\nabla \log p_t(X)] = \mathbb{E}_0 \left[\frac{1 + \langle t, g(X) \rangle}{1 + \langle t, g(X) \rangle} g(X) \right] = \mathbb{E}_0[g(X)] = 0$$

by assumption that $\mathbb{E}_0[g] = 0$. Then at $t = 0$, the Fisher information is

$$J(0) = \mathbb{E}_0[g(X)g(X)^\top] = \text{Cov}_0(g(X))$$

and

$$J(t) = \mathbb{E}_0 \left[\frac{1}{1 + \langle t, g(X) \rangle} g(X)g(X)^\top \right] = \mathbb{E}_0[g(X)g(X)^\top (1 - \langle t, g(X) \rangle)] + O(\|t\|^2) = J(0) + O(\|t\|),$$

uniformly in t near 0, because g is bounded. The information J is also continuously differentiable and positive definite as soon as $\text{Cov}_0(g(X)) \succ 0$.

As such, whenever Assumption A.12.1 holds, we have satisfied the assumptions necessary for Theorem 12.2.7. We therefore obtain the following corollary:

Corollary 12.3.2. *Let Assumption A.12.1 hold and P_t be the tilted distributions (12.3.1), and let conditions (i)–(v) hold with t replacing θ . Then for any estimator $\hat{\theta} : \mathcal{X}^n \rightarrow \mathbb{R}^d$,*

$$\begin{aligned} & \int \mathbb{E}_{P_t} \left[(\hat{\theta}(X_1^n) - \theta(P_t))^\top A(t)^{-1} (\hat{\theta}(X_1^n) - \theta(P_t)) \right] \pi(t) dt \\ & \geq \frac{\mathbb{E}_\pi [\text{tr}(\dot{\theta}(T)C(T))]^2}{n \mathbb{E}_\pi [\text{tr}(C(T)A(T)C(T)^\top J(T))] + J(\pi | A, C)}. \end{aligned}$$

As in Section 12.2, we may make particular choices of A and C to optimize the lower bound in Corollary 12.3.2. In analogy with Corollary 12.2.10, for example, if we assume the inverse Fisher information for the parameter θ ,

$$J_\theta(t) := \left(\dot{\theta}(t) J(t)^{-1} \dot{\theta}(t) \right)^{-1}$$

is continuously differentiable in t and positive definite, then we obtain

$$\mathbb{E}_\pi \left[(\hat{\theta}(X_1^n) - \theta(P_T))^\top J_\theta(T) (\hat{\theta}(X_1^n) - \theta(P_T)) \right] \geq \frac{d}{n} - O(1/n^2).$$

One particular (abstract, but common) setting is important: when $\theta(P_t)$ is not only differentiable in t , but where in fact there exists a mean-zero *influence function* $\dot{\theta}_{P_0} : \mathcal{X} \rightarrow \mathbb{R}^d$ for θ satisfying

$$\theta(P_t) = \theta(P_0) + \mathbb{E}_{P_0}[\dot{\theta}_{P_0}(X) \langle g(X), t \rangle] + o(\|t\|), \quad (12.3.2)$$

so that the derivative in Assumption A.12.1 is linear in g and $\dot{\theta}(0) = \mathbb{E}_{P_0}[\dot{\theta}_{P_0}(X)g(X)^\top]$. (The assumption that $\dot{\theta}_{P_0}$ is mean zero is no loss of generality, as $\mathbb{E}_{P_0}[g(X)] = 0$ by assumption.) This linearity holds, for example, for the mean as in Example 12.3.1, where $\dot{\theta}_{P_0}(x) = x - \mathbb{E}_{P_0}[X]$; we will also see it presently for general M-estimators in Section 12.3.1.

Let us proceed somewhat heuristically to show the types of lower bounds the existence function (12.3.2) provides. For simplicity, let us take the prior $\pi = \pi_n^d$ for a normalized prior density $\pi_n(t) = \sqrt{n}\pi_0(\sqrt{n} \cdot t)$ where π_0 is smooth and supported on $[-1, 1]$; with this we may assume the prior information $J(\pi)$ scales as $d\sqrt{n}$ (recall Corollary 12.2.2, and see also Exercise 12.3). Assuming we can vary g while maintaining the linearity (12.3.2), then, we assume $g : \mathcal{X} \rightarrow \mathbb{R}^d$ and may use $J(t) = \text{Cov}_0(g(X)) + O(\|t\|)$ and take $C = I_d$ to obtain

$$\mathbb{E}_\pi \left[(\hat{\theta}(X_1^n) - \theta(P_T))^\top A(T)^{-1} (\hat{\theta}(X_1^n) - \theta(P_T)) \right] \geq \frac{\text{tr}(\mathbb{E}_{P_0}[\dot{\theta}_{P_0}(X)g(X)^\top])^2}{n \mathbb{E}_\pi [\text{tr}(A(T)\text{Cov}_0(g(X)))]} - o(1/n)$$

for any estimator $\hat{\theta}$ and positive definite $A(t)$. The choice $A(t) = I_d$ then gives

$$\mathbb{E}_\pi \left[\|\hat{\theta}(X_1^n) - \theta(P_T)\|_2^2 \right] \geq \frac{\mathbb{E}_{P_0}[\langle \dot{\theta}_{P_0}(X), g(X) \rangle]^2}{n \mathbb{E}_{P_0}[\|g(X)\|_2^2]} - o(1/n),$$

which by Cauchy-Schwarz is maximized by $g(X) = \dot{\theta}_{P_0}(X)$, giving the lower bound

$$\mathbb{E}_\pi \left[\|\hat{\theta}(X_1^n) - \theta(P_T)\|_2^2 \right] \geq \frac{\mathbb{E}_{P_0}[\|\dot{\theta}_{P_0}(X)\|_2^2]}{n} - o(1/n). \quad (12.3.3)$$

Alternatively, if we let $\Sigma = \text{Cov}(\dot{\theta}_{P_0}(X))$, then setting $A = \Sigma$ we obtain

$$\mathbb{E}_\pi \left[\|\hat{\theta}(X_1^n) - \theta(P_T)\|_{\Sigma^{-1}}^2 \right] \geq \frac{\mathbb{E}[\langle \dot{\theta}_{P_0}(X), g(X) \rangle]^2}{n \text{tr}(\Sigma \text{Cov}_0(g(X)))} - o(1/n),$$

where $\|x\|_A = \sqrt{x^\top A x}$ is the Mahalanobis norm. Taking $g(x) = \Sigma^{-1/2} \dot{\theta}_{P_0}(x)$, we obtain

$$\mathbb{E}_\pi \left[\|\hat{\theta}(X_1^n) - \theta(P_T)\|_{\Sigma^{-1}}^2 \right] \geq \frac{d}{n} - o(1/n).$$

In comparison with the bounds available via information-theoretic techniques, such as those Exercises 9.12 and 9.13 develop, we see that we have lost some generality in that the results apply only to the squared error, but have gained in that the leading constants are unimprovable.

12.3.1 An extended example: M-estimation lower bounds

M-estimation problems provide an important class of examples to in which to apply the general nonparametric lower bounds we have developed. Here, we recall the setting of Chapter 5.3, where we have population loss $L(\theta) := \mathbb{E}_{P_0}[\ell(\theta, Z)]$ for a loss function ℓ and distribution P_0 on Z , and define the parameter

$$\theta(P_0) := \underset{\theta}{\operatorname{argmin}} L(\theta).$$

To provide lower bounds on estimating $\theta(P)$, the simplest approach is to derive an influence function expansion (12.3.2) under tilted models P_t of P_0 . To do this, we can use the implicit function theorem, which will show fairly precisely how minimizer $\theta(P_t)$ vary with tiltings (12.3.1).

Before stating the result, we define a bit of notation: let $L_t(\theta) = \mathbb{E}_{P_t}[\ell(\theta, Z)] = \mathbb{E}_0[(1 + \langle t, g(Z) \rangle) \ell(\theta, Z)]$ be the tilted population loss, and for shorthand let $\theta_t = \theta(P_t)$ for t near 0. Then applying the implicit function theorem, we have the following lemma (we defer its proof to the end of this section); the lemma holds under much weaker conditions, but we prefer to keep things simple.

Lemma 12.3.3. *Let $g : \mathcal{X} \rightarrow \mathbb{R}^k$ be a bounded and mean-zero function and Assumption A.5.1 hold on the losses ℓ and population loss L_0 . Then for all t in a neighborhood of 0, θ_t exists, is unique, and is differentiable in t , with*

$$\theta_{t+v} - \theta_t = -\mathbb{E}_0 \left[(\nabla_\theta^2 L_t(\theta_t))^{-1} \nabla \ell(\theta_t, Z) g(Z)^\top \right] v + o(\|v\|)$$

as $v \rightarrow 0$.

Immediately, the smoothness conditions in Assumption A.5.1 then imply the influence function

$$\dot{\theta}_{P_0}(z) = -(\nabla^2 L_0(\theta_0))^{-1} \nabla \ell(\theta_0, Z),$$

and moreover, we have derivatives

$$D_t \theta_t = \mathbb{E}_{P_0}[\dot{\theta}_{P_0}(Z) g(Z)^\top] + O(\|t\|).$$

The influence expansion (12.3.2) certainly holds.

Applying Corollary 12.3.2, we therefore obtain the following.

Corollary 12.3.4. *Let π_n be smooth enough and have support on the scaled ℓ_2 -ball $\frac{1}{\sqrt{n}} \mathbb{B}_2^d$ of radius $1/\sqrt{n}$. Then for any matrix $A(t) \succ 0$, continuous in t , and estimator $\hat{\theta}$,*

$$\mathbb{E}_{\pi_n} \left[\|\hat{\theta}(X_1^n) - \theta(P_t)\|_{A(t)^{-1}}^2 \right] \geq \frac{\mathbb{E}_{P_0} [\langle \nabla^2 L_0(\theta_0)^{-1} \nabla \ell(\theta_0, Z), g(Z) \rangle]^2}{n \mathbb{E}_{\pi_n} [\operatorname{tr}(A(T) \operatorname{Cov}_0(g(Z)))]} - o(1/n).$$

Exercise 12.8 shows how particular choices of g and A yield different lower bounds. In brief, however, Corollary 12.3.4 shows that asymptotically in n , the error of the empirical risk minimizers (5.3.1) is optimal.

Finally, we return to prove Lemma 12.3.3.

Proof of Lemma 12.3.3 The result is nearly immediate once we have the implicit function theorem, a version of which we state here. (We provide a proof sketch as it can be useful for simply remembering exactly what the implicit function theorem states.)

Lemma 12.3.5 (The implicit function theorem). *Let $F : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ be continuously differentiable in its coordinates $(x, y) \in \mathbb{R}^n \times \mathbb{R}^m$. Assume that $F(x_0, y_0) = 0$ and that the x -derivative matrix $D_x F(x, y) \in \mathbb{R}^{n \times n}$ is invertible at (x_0, y_0) . Then there is an open neighborhood U of y_0 such that a unique solution $x(y)$ to $F(x, y) = 0$ exists, and it is continuously differentiable on U with*

$$\dot{x}(y) = -(D_x F(x, y))^{-1} D_y F(x, y).$$

Sketch of Proof We will only show what the form of the derivative must be, presuming it exists; given the form that the derivative *should* take then an argument involving the Banach fixed point theorem demonstrates the existence. So assume that $F(x, y) = 0$; we expand the identity $0 = F(x + \Delta_x, y + \Delta_y)$ to solve for the form Δ_x must take in terms of y , and then the linear part of this form gives $\dot{x}(y)$. To that end, we have

$$\begin{aligned} 0 &= F(x + \Delta_x, y + \Delta_y) \\ &= F(x, y) + D_x F(x, y) \Delta_x + D_y F(x, y) \Delta_y + O(\|\Delta_x\|^2 + \|\Delta_y\|^2). \end{aligned}$$

Then using $F(x, y) = 0$, it must be the case that $\Delta_x = -(D_x F(x, y))^{-1} D_y F(x, y) \Delta_y + O(\|\Delta_x\|^2 + \|\Delta_y\|^2)$. Proceeding heuristically, we may take $\Delta_y \rightarrow 0$ and $\Delta_x \rightarrow 0$ to obtain that to first-order, $\Delta_x = (D_x F(x, y))^{-1} D_y F(x, y) \Delta_y + O(\|\Delta_y\|^2)$, giving $\dot{x}(y) = (D_x F(x, y))^{-1} D_y F(x, y)$. \square

To apply the implicit function theorem (Lemma 12.3.5), we compute the derivatives of $L_t(\theta) = \mathbb{E}_0[(1 + \langle t, g(Z) \rangle) \ell(\theta, Z)]$ via

$$\nabla_t L_t(\theta) = \mathbb{E}_0[g(Z) \ell(\theta, Z)], \quad \nabla_\theta L_t(\theta) = \mathbb{E}_0[(1 + \langle t, g(Z) \rangle) \nabla \ell(\theta, Z)]$$

and the second derivatives

$$\nabla_{t, \theta}^2 L_t(\theta) = \mathbb{E}_0[g(Z) \nabla \ell(\theta, Z)^\top] \in \mathbb{R}^{k \times d} \quad \text{and} \quad \nabla_\theta^2 L_t(\theta) = \mathbb{E}_0[(1 + \langle t, g(Z) \rangle) \nabla^2 \ell(\theta, Z)].$$

At $\theta_0 = \theta(P_0)$, we have by assumption that $\nabla_\theta^2 L_0(\theta) = \mathbb{E}_0[\nabla^2 \ell(\theta_0, Z)] \succ 0$, and $\nabla_\theta L_0(\theta_0) = 0$ as well. Make the identifications $x \mapsto \theta$, $y \mapsto t$ in the implicit function theorem, so we also identify $F(x, y) \mapsto \nabla_\theta L_t(\theta)$ and $D_x F(x, y) \mapsto \nabla_\theta^2 L_t(\theta)$ and $D_y F(x, y) \mapsto \nabla_{t, \theta}^2 L_t(\theta)^\top$. Then we can evidently write

$$\theta_{t+v} = \theta_t - \mathbb{E}_0 \left[(\nabla_\theta^2 L_t(\theta_t))^{-1} \nabla \ell(\theta_t, Z) g(Z)^\top \right] v + o(\|v\|)$$

for t in a neighborhood of 0 and small $v \in \mathbb{R}^k$. \square

12.4 Super-efficiency and instance optimality

Much of our development of lower bounds before this chapter was for worst-case (minimax) error measures. No statistician in their right mind should find these fully compelling—one wants to know how hard estimation and learning in the actual distribution P at hand is! Consider, for instance, the family $\mathcal{P} = \{\mathbf{N}(\theta, \sigma^2 I_d)\}_{\theta \in \mathbb{R}^d, \sigma^2 < \infty}$ of Gaussian distributions with unknown variance. Then the worst-case mean-square error for estimating the mean over this class is of course $+\infty$, while for each distribution $\mathbf{N}(\theta, \sigma^2 I_d) \in \mathcal{P}$, the sample mean satisfies $\mathbb{E}[\|\bar{X}_n - \theta\|_2^2] = \frac{d\sigma^2}{n} < \infty$. We therefore seek more nuance: an “instance-optimal” lower bound, which applies to individual probability distributions $P \in \mathcal{P}$, describing how hard an estimation problem is for the particular P at hand.

To set the stage, we highlight two desiderata that a satisfying benchmark—here, we avoid the terminology of lower bounds—of a problem’s complexity should enjoy:

- (i) The benchmark should be *instance-specific*, providing a quantity for each $P \in \mathcal{P}$
- (ii) The benchmark should be *uniformly achievable*, in that there should exist a procedure achieving the benchmark performance on *each* instance $P \in \mathcal{P}$

At some level, minimax complexity guarantees satisfy the two desiderata (i)–(ii): a constant function at least provides a quantity for each $P \in \mathcal{P}$, and minimaxity means that the result is indeed achievable. The key is that an instance-optimal bound should provide a converse, which we consider here.

- (iii) The benchmark should provide a *super-efficiency* converse: no procedure outperforms the benchmark except on a negligible collection of instances.

The bounds we have developed in the local Fano method, as in Chapter 9.4, provide bounds that at least approach the desiderata (i)–(iii): they provide a lower bound centered around a given distribution P_θ , and (frequently) these bounds are achievable. Sometimes, the bounds are independent of the parameter: Example 9.4.4 shows that in d -dimensional the normal location family $\{\mathbf{N}(\theta, \sigma^2 I_d)\}_{\theta \in \mathbb{R}^d}$, we have

$$\sup_{\|\theta - \theta_0\|_2 \leq c\sqrt{d/n}} \mathbb{E}_\theta \left[\|\hat{\theta}_n - \theta\|_2^2 \right] \gtrsim \frac{d\sigma^2}{n} \quad (12.4.1)$$

for some constant c , and the sample mean \bar{X}_n obviously achieves the lower bound. So for the squared error, the quantity $\frac{d\sigma^2}{n}$ can serve as a benchmark: it depends on the variance σ^2 of the instance at hand, and it is uniformly achievable. What is missing, however, is a super-efficiency (iii) converse: at least in the statement of the bound (12.4.1), nothing prevents an estimator from achieving much smaller than $\frac{d\sigma^2}{n}$ error at all except for a few worst-case points.

The variations on the Van Trees inequality in Sections 12.2 and 12.3 show that for parametric families (and even beyond), the Fisher information provides a super-efficiency guarantee for the squared error of the form (iii): by Corollary 12.2.9 and the bounds following, for smooth enough prior densities π , for any estimator $\hat{\theta}_n$ we have

$$\int \mathbb{E}_\theta \left[\|\hat{\theta}_n - \theta\|_2^2 \right] \pi(\theta) d\theta \geq \frac{1}{n} \int \text{tr}(J(\theta)^{-1}) \pi(\theta) d\theta - O(1/n^2),$$

and maximum-likelihood estimators attain these lower bounds (under appropriate conditions). In the case of the normal location family, because $J(\theta) = \frac{1}{\sigma^2} I_d$, we can then strengthen the lower bound (12.4.1) to

$$\inf_{\hat{\theta}_n} \int \mathbb{E}_\theta \left[\|\hat{\theta}_n - \theta\|_2^2 \right] \pi(\theta) d\theta \geq \frac{d\sigma^2}{n} (1 - o(1))$$

for “most” priors π .

To provide a bit more weight to this discussion, consider the setting of Section 12.2.3, where we wish to estimate a differentiable function $\psi(\theta) \in \mathbb{R}^p$ of the parameter $\theta \in \mathbb{R}^d$ of interest, and measure error via a quadratic with matrix $A(\theta) \succ 0$. This of course also includes the non-parametric settings in Section 12.3. For a procedure $\hat{\psi}$, or more accurately, a sequence of procedures $\hat{\psi} = \{\hat{\psi}_n\}$ defined for each sample size n , define the pointwise limiting squared error at θ by

$$L_A(\theta, \hat{\psi}) := \limsup_n n \cdot \mathbb{E}_\theta \left[(\hat{\psi}_n(X_1^n) - \psi(\theta))^\top A(\theta)^{-1} (\hat{\psi}_n(X_1^n) - \psi(\theta)) \right]$$

Theorem 12.4.1. *Let the conditions of Theorem 12.2.7 hold. Then for any estimator sequence $\hat{\psi}_n$,*

$$\int L_A(\theta) \pi(\theta) d\theta \geq \int \text{tr} \left(\dot{\psi}(\theta) J(\theta)^{-1} \dot{\psi}(\theta)^\top A(\theta)^{-1} \right) \pi(\theta) d\theta.$$

Proof In Corollary 12.2.8, define the matrix $C(\theta) = J(\theta)^{-1} \dot{\psi}(\theta)^\top A(\theta)^{-1}$, which gives for any n that

$$n\mathbb{E}_\pi \left[(\hat{\psi}_n(X_1^n) - \psi(\theta))^\top A(\theta)^{-1} (\hat{\psi}_n(X_1^n) - \psi(\theta)) \right] \geq \mathbb{E}_\pi [\text{tr}(\dot{\psi}(\theta) J(\theta)^{-1} \dot{\psi}(\theta)^\top A(\theta)^{-1})] - O(1/n),$$

where the big- O hides terms that depend on the prior π . Now define the sequence of functions

$$f_n(\theta) := n\mathbb{E}_\theta \left[(\hat{\psi}_n(X_1^n) - \psi(\theta))^\top A(\theta)^{-1} (\hat{\psi}_n(X_1^n) - \psi(\theta)) \right],$$

each of which is nonnegative. We have $\limsup_n f_n(\theta) = L_A(\theta) \in [0, \infty]$, and as $\limsup_n f_n(\theta) = \lim_{n \rightarrow \infty} \sup_{m \geq n} f_m(\theta)$, monotone convergence implies $\int \sup_{m \geq n} f_m(\theta) \pi(\theta) d\theta \downarrow \int L_A(\theta) \pi(\theta) d\theta$. \square

Taking $A = I_d$ and $\psi(\theta) = \theta$ to be the identity mapping, Theorem 12.4.1 thus shows that for “most” priors π , the limiting mean squared error of estimators satisfies

$$\int \limsup_n n\mathbb{E}_\theta \left[\|\hat{\theta}_n - \theta\|_2^2 \right] \pi(\theta) d\theta \geq \int \text{tr}(J(\theta)^{-1}) \pi(\theta) d\theta. \quad (12.4.2)$$

So, for example, there can be *no* open set U on which an estimator has mean squared error (asymptotically) better than $\frac{1}{n} \text{tr}(J(\theta)^{-1})$ everywhere on that set. Indeed, supposing to the contrary that $L(\theta) := \limsup_n n\mathbb{E}_\theta [\|\hat{\theta}_n - \theta\|_2^2] < \text{tr}(J(\theta)^{-1})$ for each $\theta \in U$, we take any smooth prior whose support U contains, and then obtain the contradiction that

$$\int L(\theta) \pi(\theta) d\theta < \int \text{tr}(J(\theta)^{-1}) \pi(\theta) d\theta.$$

Another perspective on this is that no “good” estimator can outperform. Suppose that an estimator is efficient in some neighborhood of θ_0 for the squared error, meaning that

$$\limsup_n n\mathbb{E}_\theta \left[\|\hat{\theta}_n - \theta\|_2^2 \right] \leq \text{tr}(J(\theta)^{-1})$$

for each θ near θ_0 . Then the set of points θ at which strict inequality can hold above necessarily has measure 0: rearranging the preceding inequality we have

$$0 \stackrel{(\star)}{\geq} \int \left(\text{tr}(J(\theta)^{-1}) - \limsup_n n \mathbb{E}_\theta \left[\|\hat{\theta}_n - \theta\|_2^2 \right] \right) \pi(\theta) d\theta \geq 0$$

where the inequality (\star) is inequality (12.4.2). As the integrand is nonnegative, it must be 0 for almost every θ . Such average-case super-efficiency converses are also present in the classical theory of asymptotic estimator efficiency that Le Cam and Hajak develop [cf. 184], which extends the present results to general loss functions far beyond the squared error.

JCD Comment: We should maybe have a small section here on super-efficiency and estimation. Perhaps more in exercises as well, assuming we do Le Cam in exercises

12.5 Applications in privacy

JCD Comment: Develop score attacks as in Cai et al. [45], Cai et al. [46].

12.6 Bibliography and further reading

JCD Comment: Discuss MMSE estimators, local asymptotic normality, and so on. Provide a few exercises on local asymptotic normality as well if possible.

12.7 Exercises

Exercise 12.1: Let Θ be a compact convex set with non-empty interior (a convex body), and let $\{P_\theta\}_{\theta \in \Theta}$ be a family of distributions, all absolutely continuous with respect to one another.

- (a) Show that if $\hat{\theta}$ is an unbiased estimator of θ , then for each $\theta \in \Theta$, $P_\theta(\hat{\theta}(X) \notin \Theta) > 0$. *Hint.* Extend Example 12.1.2 to show that if $\hat{\theta}$ is unbiased and $\hat{\theta} \in \Theta$, then Θ cannot have an interior. It may be useful to leverage supporting hyperplanes of convex bodies.
- (b) Show that if $\hat{\theta}$ is an estimator with $P_\theta(\hat{\theta}(X) \notin \Theta) > 0$, then the projected estimator $\tilde{\theta}(x) = \text{Proj}_\Theta(\hat{\theta}(x))$ satisfies

$$\mathbb{E}_\theta \left[\|\tilde{\theta}(X) - \theta\|_2^2 \right] < \mathbb{E}_\theta \left[\|\hat{\theta}(X) - \theta\|_2^2 \right]$$

Hint. Theorem B.1.11 in Appendix B.1.2 provides a useful characterization of the projection onto a convex set.

- (c) Conclude that if Θ is a convex body and $\{P_\theta\}_{\theta \in \Theta}$ satisfies the assumed conditions of the exercise, then the projected estimator $\tilde{\theta}$ satisfies

$$\mathbb{E}_\theta \left[\|\tilde{\theta}(X) - \theta\|_2^2 \right] < \mathbb{E}_\theta \left[\|\hat{\theta}(X) - \theta\|_2^2 \right]$$

for all $\theta \in \Theta$.

JCD Comment: Do a computational version of the James-Stein estimator perhaps after its discussion above.

Exercise 12.2: Let $Y = X\theta + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, I_n)$ and $X \in \mathbb{R}^{n \times d}$. For $\lambda > 0$ define $\hat{\theta}_\lambda = \operatorname{argmin}_\theta \{\|Xt - Y\|_2^2 + \lambda \|t\|_2^2\}$. Define the risk $R(\lambda) := \mathbb{E}[\|\hat{\theta}_\lambda - \theta\|_2^2]$ and $R(0) = \lim_{\lambda \downarrow 0} R(\lambda)$. Give

$$R'(0^+) := \lim_{\lambda \downarrow 0} \frac{R(\lambda) - R(0)}{\lambda},$$

show that $R'(0^+) < 0$, and conclude that for any θ , there exists $\lambda > 0$ such that $R(\lambda) < R(0)$.

Exercise 12.3: In this question, you investigate some examples and properties of the information of compactly supported \mathcal{C}^∞ functions.

- (a) Let $\phi(t) = \exp(-\frac{1}{1-t^2})$ for $|t| < 1$ and $\phi(t) = 0$ for $|t| \geq 1$. Show that ϕ is \mathcal{C}^∞ .
- (b) Show that for ϕ as in part (a), $\int_{-1}^1 \phi'(t)^2 / \phi(t) dt < \infty$.
- (c) Let ϕ be a symmetric, nonnegative, \mathcal{C}^∞ function with support $[-1, 1]$ and $\int \phi(t) dt = 1$. For $0 < a \leq 1$ define $\pi_a(t) = \phi(t/a)/a$. Show that π is a density on $[-a, a]$ and give $J(\pi_a)$ as a function of $J(\pi_1)$.

Exercise 12.4: Prove Corollary 12.2.4.

Exercise 12.5: Demonstrate the claim of Example 12.2.5 that

$$t^* := \operatorname{argmin}_t \inf_{\beta_0, \beta} \mathbb{E} \left[(Y - \beta_0 - tZ - \beta^\top X)^2 \right]$$

satisfies $t^* = \mathbb{E}[Y(1)] - \mathbb{E}[Y(0)]$.

Exercise 12.6: In conditions (i)–(v) for Theorem 12.2.7, replace all instances of “suitably continuous” with “continuously differentiable,” assume that $J(\theta) = \mathbb{E}_\theta[\dot{\ell}_\theta(X)\dot{\ell}_\theta(X)^\top]$ is continuous in θ , and assume instead of condition (v) that π is a product measure on $[a_1, b_1] \times \cdots \times [a_d, b_d]$, as in Theorem 12.2.6. Prove Theorem 12.2.7 under these conditions.

Exercise 12.7 (Estimation with and without nuisance parameters): Consider probabilistic models $\{P_{\theta, \eta}\}$, $\theta \in \mathbb{R}^d$ and $\eta \in \mathbb{R}^k$, with densities $p_{\theta, \eta}$ and (joint) Fisher information matrix

$$J(\theta, \eta) = \begin{bmatrix} J_{\theta, \theta} & J_{\theta, \eta} \\ J_{\theta, \eta}^\top & J_{\eta, \eta} \end{bmatrix},$$

where the notation indicates partitioning J to match parameters.

- (a) Define the block matrices

$$M := \begin{bmatrix} A & B \\ B^\top & C \end{bmatrix} \quad \text{and} \quad M^{-1} = \begin{bmatrix} X & Y \\ Y^\top & Z \end{bmatrix},$$

where the matrices are partitioned similarly. Assume that $M \succ 0$. Show that $X \preceq A^{-1}$, with equality if and only if $B = 0$.

(b) Show that

$$J_{\theta,\theta}^{-1} \preceq [J(\theta, \eta)^{-1}]_{\theta,\theta}.$$

When does equality hold?

(c) Use Corollary 12.2.10 to interpret the above inequality in the context of problems in which the nuisance parameters η are known versus unknown.

Exercise 12.8: Consider M-estimation problems as in Corollary 12.3.4. Let $H = \nabla^2 L(\theta_0)$ be the Hessian of L at $\theta_0 = \operatorname{argmin}_{\theta} L(\theta)$, and let $\Sigma = \operatorname{Cov}(\nabla \ell(\theta_0, Z))$ the covariance of the gradients. Give appropriate choices of the matrices $A(t)$ and tilting functions $g : \mathcal{Z} \rightarrow \mathbb{R}^d$ to show the following:

(a) For smooth enough priors π_n on T supported on $\frac{1}{\sqrt{n}}\mathbb{B}_2^d$, any estimator $\hat{\theta}$ satisfies

$$\liminf_n n \cdot \mathbb{E}_{\pi_n} \left[\|\hat{\theta}(Z_1^n) - \theta(P_T)\|^2 \right] \geq \operatorname{tr}(H^{-1}\Sigma H^{-1}).$$

(b) For smooth enough priors π_n on T supported on $\frac{1}{\sqrt{n}}\mathbb{B}_2^d$, any estimator $\hat{\theta}$ satisfies

$$\liminf_n n \cdot \mathbb{E}_{\pi_n} \left[(\hat{\theta}(Z_1^n) - \theta(P_T))^\top H \Sigma^{-1} H (\hat{\theta}(Z_1^n) - \theta(P_T)) \right] \geq d.$$

What do these results say about the empirical risk minimizer (5.3.1)?

Chapter 13

Testing and functional estimation

When we wish to estimate a complete “object,” such as the parameter θ in a linear regression $Y = X\theta + \varepsilon$, or a density when we observe X_1, \dots, X_n i.i.d. with a density f , the previous chapters give a number of approaches to proving fundamental optimality results and limits. In many cases, however, we wish to estimate *functionals* of a distribution or larger parameter, rather than the entire distribution or a high-dimensional parameter. Suppose we wish to estimate some statistic $T(P) \in \mathbb{R}$ of a probability distribution P . Then a naive estimator is to construct an estimate \hat{P} of P , and simply plug it in: use $\hat{T} = T(\hat{P})$. But frequently—and as we have seen in the preceding chapters—our ability to estimate \hat{P} may be limited, while various statistics of P may be easier to estimate. As a trivial example of this phenomenon, suppose we have an unknown distribution P supported on $[-1, 1]$, and we wish to estimate the statistic $T(P) = \mathbb{E}_P[X]$, its expectation. Then the trivial sample mean estimator

$$T_n := \bar{X}_n$$

satisfies $\mathbb{E}[(T_n - \mathbb{E}[X])^2] \leq \frac{1}{n}$. But an estimator that first attempts to approximate the full distribution P via some \hat{P} and then estimate $\int x d\hat{P}(x)$ is likely to incur substantial additional error.

Alternatively, we might wish to test different properties of distributions. In *goodness of fit testing*, we are given a sample X_1, \dots, X_n i.i.d. from a distribution Q , and we wish to distinguish whether $Q = P$ or Q is far from P . In related *two-sample tests*, we are given samples $X_1^n \stackrel{\text{iid}}{\sim} P$ and $Y_1^m \stackrel{\text{iid}}{\sim} Q$, and again wish to test whether $Q = P$ or Q and P are far from one another. For example, in a medical study, we may wish to distinguish whether there are significant differences between a treated population Q and control population P .

More broadly, we wish to develop tools to understand the optimality of different estimators and tests of *functionals*, by which we mean scalar valued parameters of a distribution P . Such parameters could include the norm $\|\theta\|_2$ of a regression vector, an estimate of the best possible expected loss $\inf_f \mathbb{E}_P[\ell(f(X), Y)]$ in a prediction problem, the distance $\|P - P_0\|_{\text{TV}}$ of a sampled population P from a reference P_0 , or the probability mass of outcomes we have not observed in a study. This chapter develops a few of the tools to understand these problems.

13.1 Geometrizing rates of convergence

JCD Comment: Figure on modulus of continuity?

In some cases, it is possible to reduce the development of lower bounds for estimation problems to characterizing purely geometric objects, relating the continuity of a function to be estimated to distances of the underlying probability distributions. This approach makes the intuition that estimation should be hard when the statistic of interest is quite sensitive to the underlying distribution quantitative. To proceed, we define the *local modulus of continuity* of a functional $\theta : \mathcal{P} \rightarrow \mathbb{R}$ on a family \mathcal{P} of distributions at a fixed distribution P_0 with respect to the Hellinger distance,

$$\omega_{\text{hel}}(\epsilon; \theta, P_0, \mathcal{P}) := \sup_{P_1 \in \mathcal{P}} \{|\theta(P_0) - \theta(P_1)| \mid d_{\text{hel}}(P_0, P_1) \leq \epsilon\}. \quad (13.1.1)$$

We recognize this as the local Lipschitz constant of the parameter θ for the Hellinger distance around a fixed distribution P_0 . The *global modulus of continuity* is simply the supremum of the local modulus over all $P_0 \in \mathcal{P}$,

$$\omega_{\text{hel}}(\epsilon; \theta, \mathcal{P}) := \sup_{P \in \mathcal{P}} \omega_{\text{hel}}(\epsilon; \theta, P, \mathcal{P}).$$

These are “geometric” quantities in that each measures how much θ can vary over small neighborhoods with respect to the Hellinger distance. Using Le Cam’s two point method (Chapter 9.3), we can nearly immediately see that this Hellinger modulus implies lower bounds on estimation error. In this case, we lower bound a somewhat smaller quantity than the minimax error, instead bounding the “hardest one-dimensional sub-problem,”

$$\mathfrak{M}_n^{\text{1d}}(\theta(\mathcal{P}), P_0, \Phi) := \sup_{P_1 \in \mathcal{P}} \inf_{\hat{\theta}_n} \max_{P \in \{P_0, P_1\}} \mathbb{E}_{P^n} \left[\Phi(|\hat{\theta}_n(X_1^n) - \theta(P)|) \right]. \quad (13.1.2)$$

The quantity (13.1.2) differs from the standard minimax risk in that the distribution P_0 is fixed, and nature must first choose an alternative P_1 ; the estimator $\hat{\theta}_n$ then *knows* P_1 and P_0 . Any lower bound on the hardest subproblem (13.1.2) immediately lower bounds the minimax risk, as

$$\mathfrak{M}_n(\theta(\mathcal{P}), \Phi) \geq \sup_{P_0 \in \mathcal{P}} \mathfrak{M}_n^{\text{1d}}(\theta(\mathcal{P}), P_0, \Phi).$$

Proposition 13.1.1. *Let $\Phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ be any nondecreasing loss function and $\theta : \mathcal{P} \rightarrow \mathbb{R}$. Define the sequence $\epsilon_n^2 = \frac{2-\sqrt{2}}{n}$. Then for $n \geq 2$,*

$$\mathfrak{M}_n^{\text{1d}}(\theta(\mathcal{P}), P_0, \Phi) \geq \frac{1}{2} \left(1 - \sqrt{3/8}\right) \Phi \left(\frac{1}{2} \omega_{\text{hel}}(\epsilon_n; \theta, P_0, \mathcal{P}) \right).$$

Additionally, for all $c > 0$,

$$\mathfrak{M}_n^{\text{1d}}(\theta(\mathcal{P}), P_0, \Phi) \geq \frac{1 - \sqrt{1 - e^{-c}} \sqrt{1 - e^{-2c}}}{2} \cdot \Phi \left(\frac{1}{2} \omega_{\text{hel}} \left(\sqrt{\frac{c}{n}}; \theta, P_0, \mathcal{P} \right) \right) \cdot (1 - o(1)).$$

Proof This is nearly immediate from the minimax lower bound (9.3.3) from Le Cam’s two point method, which implies

$$\mathfrak{M}_n^{\text{1d}}(\theta(\mathcal{P}), P_0, \Phi) \geq \frac{1}{2} \sup_{P_1 \in \mathcal{P}} \Phi(|\theta(P_0) - \theta(P_1)|/2) (1 - \|P_0^n - P_1^n\|_{\text{TV}}).$$

Now we use Proposition 2.2.7 to observe that

$$\begin{aligned} \|P_0^n - P_1^n\|_{\text{TV}} &\leq d_{\text{hel}}(P_0^n, P_1^n) \sqrt{2 - d_{\text{hel}}^2(P_0^n, P_1^n)} \\ &= (1 - (1 - d_{\text{hel}}^2(P_0, P_1))^n) \sqrt{1 + (1 - d_{\text{hel}}^2(P_0, P_1))^n} \end{aligned}$$

by the tensorization identity (9.2.4) for Hellinger distance. Because $(1 - c/n)^n \rightarrow e^{-c}$, the asymptotic limit in the proposition follows from the inequality

$$\limsup_n \sup_{d_{\text{hel}}^2(P, Q) \leq \frac{c}{n}} \|P^n - Q^n\|_{\text{TV}} \leq (1 - e^{-c}) \sqrt{1 + e^{-c}} = \sqrt{1 - e^{-c}} \sqrt{1 - e^{-2c}}.$$

Let us perform the evaluation in finite samples to obtain explicit constants. Because $x \mapsto x\sqrt{2-x}$ is increasing for $0 \leq x \leq 1$, we see that if $d_{\text{hel}}^2(P_0, P_1) \leq \frac{c}{n}$ then

$$\|P_0^n - P_1^n\|_{\text{TV}} \leq \left(1 - \left(1 - \frac{c}{n}\right)^n\right) \sqrt{1 + \left(1 - \frac{c}{n}\right)^n}.$$

If $(1 - \frac{c}{n})^n \geq \frac{1}{2}$, then $\|P_0^n - P_1^n\|_{\text{TV}} \leq \frac{1}{2} \sqrt{3/2} = \sqrt{3/8}$. Because $(1 - \frac{c}{n})^n \uparrow e^{-c}$ as n grows, taking $n = 2$ we solve $(1 - c/2)^2 = \frac{1}{2}$, or $c = 2 - \sqrt{2}$, to obtain that $d_{\text{hel}}^2(P_0, P_1) \leq \frac{2-\sqrt{2}}{n}$ implies $\|P_0^n - P_1^n\|_{\text{TV}} \leq \sqrt{3/8}$ for $n \geq 2$. Substituting into the minimax lower bound gives the result. \square

So long as the local modulus does not vary too wildly as $\epsilon \rightarrow 0$, it in fact *characterizes* the difficulty of the hardest one-dimensional subproblem at P_0 . We say that the modulus is regular at P_0 if there exist $0 < r_1 \leq r_0$ and K_0, K_1 such that

$$K_0 \epsilon^{r_0} \leq \omega_{\text{hel}}(\epsilon; \theta, P_0, \mathcal{P}) \leq K_1 \epsilon^{r_1}$$

for all small $\epsilon > 0$. Then we have the following complement to Proposition 13.1.1, showing that the Hellinger modulus indeed is the fundamental quantity governing the risk (13.1.2) of the hardest one-dimensional subproblem.

Proposition 13.1.2. *The risk of the hardest one-dimensional subproblem satisfies*

$$\mathfrak{M}_n^{\text{ld}}(\theta(\mathcal{P}), P_0, \Phi) \leq \sup_{\epsilon \geq 0} e^{-n\epsilon^2} \Phi(\omega(\epsilon)). \quad (13.1.3)$$

If additionally $\Phi(t) = t^p$ for some $p > 0$ and the local modulus is regular at P_0 , then for large enough n ,

$$\mathfrak{M}_n^{\text{ld}}(\theta(\mathcal{P}), P_0, \Phi) \leq e^{-\frac{1}{2}r_1 p} \cdot \Phi\left(\omega_{\text{hel}}\left(\sqrt{\frac{r_0 p}{2n}}; \theta, P_0, \mathcal{P}\right)\right).$$

Proof Let $\omega(\epsilon) = \omega_{\text{hel}}(\epsilon; \theta, P_0, \mathcal{P})$ for shorthand. Given P_0, P_1 , let $\theta_0 = \theta(P_0)$ and $\theta_1 = \theta(P_1)$, and let Ψ_n be the optimal test between P_0^n and P_1^n , so that

$$P_0(\Psi_n \neq 0) + P_1(\Psi_n \neq 1) = 1 - \|P_0^n - P_1^n\|_{\text{TV}} \leq 1 - d_{\text{hel}}^2(P_0^n, P_1^n) = (1 - d_{\text{hel}}^2(P_0, P_1))^n$$

by Proposition 2.2.7 and the tensorization identity (9.2.4). Define the estimator $\hat{\theta}_n = \theta_1$ if $\Psi_n = 1$ and $\hat{\theta}_n = \theta_0$ otherwise. Then

$$\max_{P \in \{P_0, P_1\}} \mathbb{E}_P \left[\Phi(|\hat{\theta}_n - \theta(P)|) \right] \leq (1 - d_{\text{hel}}^2(P_0, P_1))^n \Phi(|\theta_0 - \theta_1|).$$

Because $(1 - c)^n \leq e^{-nc}$, we therefore obtain inequality (13.1.3).

By the regularity assumption on ω_{hel} , the ϵ_n^* minimizing the upper bound (13.1.3) satisfy

$$\sqrt{\frac{r_1 p}{2n}} \leq \epsilon_n^* \leq \sqrt{\frac{r_0 p}{2n}}$$

(take the ϵ values minimizing, respectively, $e^{-n\epsilon^2} K_1^p \epsilon^{r_1 p}$ and $e^{-n\epsilon^2} K_0^p \epsilon^{r_0 p}$). Substituting gives the result. \square

Noting that $\frac{1}{2} - \frac{1}{2}\sqrt{1 - e^{-1}}\sqrt{1 - e^{-2}} > 1/8$, Proposition 13.1.1 implies that for large enough n , we have the (global) minimax lower bound

$$\mathfrak{M}_n(\theta(\mathcal{P}), \Phi) \geq \frac{1}{8} \cdot \Phi \left(\frac{1}{2} \omega_{\text{hel}} \left(\frac{1}{\sqrt{n}}; \theta, \mathcal{P} \right) \right). \quad (13.1.4)$$

So the modulus of continuity of the parameter θ at a radius of roughly $\frac{1}{\sqrt{n}}$ for the Hellinger provides a lower bound on estimation, reducing the problem of obtaining a minimax lower bound to one of lower bounding the modulus of continuity of θ . (The question of attainability of the implied bound is more involved; see the bibliographic section for some discussion of this and related issues.)

Example 13.1.3 (Estimating the value of a nonparametric function): Let us revisit the nonparametric regression problems of Section 10.1. Assume we receive observations $Y_i = f(X_i) + \varepsilon_i$, where ε_i are i.i.d. $\mathbf{N}(0, 1)$, $X_i \stackrel{\text{iid}}{\sim} \text{Uniform}([-1, 1])$, and f is a 1-Lipschitz function. We wish to estimate the value of f at 0, i.e., $\theta(P) = f(0)$. Once we provide a lower bound on the Hellinger modulus in this observation model, Proposition 13.1.1 then gives a minimax lower bound.

We now construct a particular function f . Let $\phi(x) = [1 - |x|]_+$, which is 1-Lipschitz, and for $t \in [0, 1]$ define $f_t(x) = t\phi(x/t)$, which is 1-Lipschitz (and f_0 is identically 0). Letting P_t and P_0 denote the joint distributions of (X, Y) when $f = f_t$ or $f = f_0$, we have

$$d_{\text{hel}}^2(P_t, P_0) = \frac{1}{2} \int_{-1}^1 d_{\text{hel}}^2(\mathbf{N}(f_t(x), 1), \mathbf{N}(f_0(x), 1)) dx.$$

The Hellinger distance between two Gaussians (see Exercise 2.2) satisfies

$$d_{\text{hel}}^2(\mathbf{N}(\mu_0, \sigma^2), \mathbf{N}(\mu_1, \sigma^2)) = 1 - \exp \left(-\frac{1}{8\sigma^2} (\mu_0 - \mu_1)^2 \right) \leq \frac{(\mu_0 - \mu_1)^2}{8\sigma^2},$$

where we use that $e^x \geq 1 + x$, or $1 - e^x \leq -x$ for all $x \in \mathbb{R}$. We therefore obtain

$$d_{\text{hel}}^2(P_t, P_0) \leq \frac{1}{16} \int_{-1}^1 f_t(x)^2 dx = \frac{1}{16} \int_{-t}^t t^2 \phi^2(x/t) dx = \frac{t^3}{16} \int_{-1}^1 \phi^2(x) dx = \frac{t^3}{24}.$$

(The factor 24 is, of course, unimportant.) Observing that the separation between the parameters of interest is $\theta(P_t) - \theta(P_0) = t$, the modulus has lower bound

$$\omega_{\text{hel}}(\epsilon; \theta, \mathcal{P}) \geq \sup \{t \mid t^3 \leq 24\epsilon\} = \sqrt[3]{24} \cdot \epsilon^{1/3}.$$

Substituting in Proposition 13.1.1 with $\epsilon_n = \frac{1}{\sqrt{n}}$, we see obtain a minimax error bound in squared error of

$$\mathfrak{M}_n(\theta(\mathcal{P}), (\cdot)^2) \gtrsim \frac{1}{n^{2/3}}$$

for the class of Lipschitz functions. Estimators (such as the kernel smoothing estimators of Chapter 10.1) achieve this rate, showing that it is sharp. \diamond

13.1.1 Fisher information and divergence measures

We have seen the Fisher information matrix appear in lower bounds in previous chapters, and at least for parametric models, we can frequently characterize the local modulus of continuity using it. This perspective will also hold beyond moduli with respect to Hellinger distance, so that for suitably regular families of distributions, f -divergences with twice differentiable f will all be locally equivalent.

The key results will show quadratic expansions of Hellinger and other divergence measures under appropriate regularity conditions on the family $\mathcal{P} = \{P_\theta\}$ of distributions. To motivate things, we begin heuristically, providing the rigorous regularity conditions presently. Assume that the distributions $P_\theta \in \mathcal{P}$ have densities p_θ , and recall the Fisher score $\dot{\ell}_\theta := \nabla_\theta \log p_\theta$, which is (typically) mean-zero. Let f be any twice differentiable convex function with $f(1) = 0$ and $f''(1) > 0$. Then for v small, we have $p_{\theta+v}(x)/p_\theta(x) - 1 = o_x(\|v\|)$, where $o_x(\|v\|)$ means that a term that may depend on x but satisfies $o_x(\|v\|)/\|v\| \rightarrow 0$ as $v \rightarrow 0$. Then

$$f\left(\frac{p_{\theta+v}(x)}{p_\theta(x)}\right) = f(1) + f'(1)\left(\frac{p_{\theta+v}(x)}{p_\theta(x)} - 1\right) + \frac{f''(1)}{2}\left(\frac{p_{\theta+v}(x)}{p_\theta(x)} - 1\right)^2 + o_x(\|v\|^2).$$

Assuming we may ignore the remainder terms (which will require some type of domination to actually be able to integrate them), we obtain

$$\begin{aligned} D_f(P_{\theta+v}\|P_\theta) &= \int f\left(\frac{p_{\theta+v}(x)}{p_\theta(x)}\right) p_\theta(x) dx \\ &= \int f'(1)(p_{\theta+v}(x) - p_\theta(x)) dx + \frac{f''(1)}{2} \int \left(\frac{p_{\theta+v}(x)}{p_\theta(x)} - 1\right)^2 dx + \int o_x(\|v\|^2) p_\theta(x) dx \\ &\stackrel{(?)}{=} \frac{f''(1)}{2} D_{\chi^2}(P_{\theta+v}\|P_\theta) + o(\|v\|^2), \end{aligned} \tag{13.1.5}$$

where the equality (?) is heuristic in that it integrates the remainder. Continuing with our heuristics, we use the expansion

$$\frac{p_{\theta+v}(x)}{p_\theta(x)} - 1 = \frac{\nabla_\theta p_\theta(x)^\top v + o_x(\|v\|)}{p_\theta(x)} = \dot{\ell}_\theta(x)^\top v + o_x(\|v\|)$$

to obtain

$$\begin{aligned} D_{\chi^2}(P_{\theta+v}\|P_\theta) &= \int \left(\frac{p_{\theta+v}(x)}{p_\theta(x)} - 1\right)^2 p_\theta(x) dx \\ &= \int \left(\dot{\ell}_\theta(x)^\top v + o_x(\|v\|)\right)^2 p_\theta(x) dx \\ &\stackrel{(?)}{=} \int (\dot{\ell}_\theta(x)^\top v)^2 p_\theta(x) dx + o(\|v\|^2) = v^\top J(\theta) v + o(\|v\|^2), \end{aligned} \tag{13.1.6}$$

where once again equality (?) is heuristic. Expressions (13.1.5) and (13.1.6) show that, at least under some regularity conditions, we expect that

$$D_f(P_{\theta+v}\|P_\theta) = \frac{f''(1)}{2} v^\top J(\theta) v + o(\|v\|^2), \tag{13.1.7}$$

so that the (local) geometry on probability distributions the Fisher information matrix $J(\theta)$ induces is equivalent to that for f -divergences.

Under the condition (13.1.7), it becomes rather straightforward to compute the modulus of continuity. Indeed, we may generalize the Hellinger modulus to define

$$\omega_f(\epsilon; T, P_\theta, \mathcal{P}) := \sup \{ |T(\theta') - T(\theta)| \mid D_f(P_{\theta'} \| P_\theta) \leq \epsilon^2 \},$$

where we square ϵ to match the Hellinger case (13.1.1), which corresponds to $f(t) = \frac{1}{2}(\sqrt{t} - 1)^2$. Then the Fisher information (asymptotically) characterizes the modulus of continuity whenever we have the identifiability condition that for all $\epsilon > 0$,

$$\inf_{\theta'} \{ D_f(P_{\theta'} \| P_\theta) \mid \|\theta - \theta'\|_2 \geq \epsilon \} > 0. \quad (13.1.8)$$

Proposition 13.1.4. *Let $\{P_\theta\}$ be a family of distributions for which the f -divergence satisfies the expansion (13.1.7) at $\theta \in \mathbb{R}^d$, and let $T : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable at θ . Then*

$$\omega_f(\epsilon; T, P_\theta, \mathcal{P}) \geq \sqrt{2/f''(1)} \cdot \epsilon \left\| J(\theta)^{-1/2} \nabla T(\theta) \right\|_2 - o(\epsilon).$$

If in addition the identifiability condition (13.1.8) holds, then

$$\omega_f(\epsilon; T, P_\theta, \mathcal{P}) \leq \sqrt{2/f''(1)} \cdot \epsilon \left\| J(\theta)^{-1/2} \nabla T(\theta) \right\|_2 + o(\epsilon).$$

Proof Given the expansion (13.1.7), as $\epsilon \downarrow 0$, for any $v \in \mathbb{R}^d$ we have

$$\frac{1}{\epsilon^2} D_f(P_{\theta+\epsilon v} \| P_\theta) \rightarrow \frac{f''(1)}{2} v^\top J(\theta) v.$$

Then because $T(\theta + v) = T(\theta) + \langle \nabla T(\theta), v \rangle + o(\|v\|)$, we have

$$\liminf_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \omega(\epsilon) \geq \sup \left\{ \langle \nabla T(\theta), v \rangle \mid v^\top J(\theta) v \leq \frac{2}{f''(1)} \right\} = \sqrt{2/f''(1)} \left\| J(\theta)^{-1/2} \nabla T(\theta) \right\|_2.$$

For the upper bound, by the identifiability assumption (13.1.8), there exists a function $\delta(\epsilon) \rightarrow 0$ as $\epsilon \downarrow 0$ for which $D_f(P_{\theta'} \| P_\theta) \leq \epsilon^2$ implies that $\|\theta' - \theta\|_2 \leq \delta(\epsilon)$. The identity (13.1.7) shows that for small enough $\epsilon > 0$, we have $D_f(P_{\theta'} \| P_\theta) = \frac{f''(1)}{2} (\theta' - \theta)^\top J(\theta) (\theta' - \theta) + o(\|\theta' - \theta\|^2)$ whenever $\|\theta' - \theta\| \leq \delta(\epsilon)$, so that in fact it must be the case that $D_f(P_{\theta'} \| P_\theta) \leq \epsilon^2$ implies that $\theta' = \theta + v$ for some v satisfying $v^\top J(\theta) v \leq \frac{2}{f''(1)} \epsilon^2 (1 + o(1))$ as $\epsilon \rightarrow 0$. So

$$\begin{aligned} \omega(\epsilon) &\leq \sup \left\{ |T(\theta + v) - T(\theta)| \mid v^\top J(\theta) v \leq (2/f''(1)) \epsilon^2 (1 + o(1)) \right\} \\ &= \sup \left\{ |\langle \nabla T(\theta), v \rangle| + o(\|v\|) \mid v^\top J(\theta) v \leq (2/f''(1)) \epsilon^2 (1 + o(1)) \right\} \\ &\leq \sup \left\{ \left\| J(\theta)^{-1/2} \nabla T(\theta) \right\|_2 \left\| J(\theta)^{1/2} v \right\|_2 + o(\|v\|) \mid \left\| J(\theta)^{1/2} v \right\|_2^2 \leq (2/f''(1)) \epsilon^2 (1 + o(1)) \right\} \end{aligned}$$

by Cauchy-Schwarz, giving the result. \square

In a sense, then, so long as we have enough regularity, most divergence measures and related moduli of continuity induce the same geometry via the Fisher information. The next subsection develops conditions under which these Fisher information expansions hold, but here we preview them a bit by claiming that the expansion (13.1.7) holds for the squared Hellinger distance, corresponding to $f(t) = \frac{1}{2}(\sqrt{t} - 1)^2 = \frac{t}{2} - \sqrt{t} + 1$ for most distributions, where $f''(1) = \frac{1}{4}$.

Corollary 13.1.5. *Let $\{P_\theta\}_{\theta \in \mathbb{R}}$ be a suitably regular family of distributions, and let $J(\theta) = \mathbb{E}_\theta[\dot{\ell}_\theta^2]$ be the Fisher information. Then there exists a numerical constant $c > 0$ such that for any θ_0 ,*

$$\sup_{\theta_1} \inf_{\hat{\theta}_n} \max_{\theta \in \{\theta_0, \theta_1\}} \mathbb{E}_\theta \left[|\hat{\theta}_n(X_1^n) - \theta| \right] \geq c \frac{1}{\sqrt{nJ(\theta_0)}}$$

for all large enough n .

13.1.2 Valid asymptotic information expansions of divergences

In this subsection, we collect a few representative regularity conditions on the distribution family to make the expansions of f -divergences in terms of the Fisher information matrix, as in the heuristic equality (13.1.7), rigorous. For the Hellinger distance, we can provide a fairly general result, which allows non-differentiable densities (so, for example, we may consider the Laplace distribution with density $p_\theta(x) = \frac{1}{2} \exp(-|x - \theta|)$).

Lemma 13.1.6. *Let $\{P_\theta\}_{\theta \in \Theta}$ be a parametric family where each P_θ has a density p_θ w.r.t. a base measure μ . Assume that $\sqrt{p_\theta(x)}$ is continuously differentiable in a neighborhood of θ_0 for μ -almost all x (though this neighborhood may depend on x). Assume as well that the Fisher information matrix $J(\theta) = \mathbb{E}_\theta[\dot{\ell}_\theta \dot{\ell}_\theta^\top]$ is continuous in θ at θ_0 . Then*

$$d_{\text{hel}}^2(P_{\theta_0+v}, P_{\theta_0}) = \frac{1}{8} v^\top J(\theta_0) v + o(\|v\|^2) \quad \text{uniformly in } v \text{ near } 0.$$

Because the measure-theoretic details are not our main focus, we defer the proof to Section 13.5.1.

The Hellinger distance admits particularly unrestrictive conditions on the family $\{P_\theta\}$, because the associated f -divergence, with $f(t) = \frac{1}{2}(\sqrt{t} - 1)^2$, grows only linearly as $t \rightarrow \infty$. When the divergences can grow more quickly, we will use stronger regularity conditions on the model family.

Definition 13.1. *Let $\{P_\theta\}_{\theta \in \Theta}$, where $\Theta \subset \mathbb{R}^d$, be a family of distributions. The family is sufficiently regular for the χ^2 -divergence at the point $\theta_0 \in \text{int } \Theta$ if the following hold:*

- i. *There exists a base measure μ for which P_θ has density p_θ for all θ in a neighborhood of θ_0 .*
- ii. *The Fisher information $J_{\theta_0} := \mathbb{E}_{\theta_0}[\dot{\ell}_{\theta_0} \dot{\ell}_{\theta_0}^\top]$ exists.*
- iii. *There is some $h : \mathcal{X} \rightarrow [0, \infty]$ with $\mathbb{E}_{\theta_0}[h^2] < \infty$ for which $|\frac{p_{\theta_0+v}(x)}{p_{\theta_0}(x)} - 1 - \dot{\ell}_{\theta_0}(x)^\top v| \leq h(x) \|v\|$ for all v in a neighborhood of 0.*

Many common distribution families are sufficiently regular.

Example 13.1.7: Consider the collection $\{N(\theta, 1)\}$ of normal distributions with densities $p_\theta(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}(x - \theta)^2)$. Then the score function is $\dot{\ell}_\theta(x) = (x - \theta)$, and (w.l.o.g.) taking $\theta = 0$, we have

$$\frac{p_v(x)}{p_0(x)} - 1 - xv = e^{-\frac{1}{2}v^2} e^{xv} - 1 - xv = (e^{-\frac{1}{2}v^2} - 1)(1 + xv) + e^{-\frac{1}{2}v^2} \sum_{k=2}^{\infty} \frac{(xv)^k}{k!}.$$

As $|e^{-\frac{1}{2}v^2} - 1| \leq v^2$ and

$$\sum_{k=2}^{\infty} \frac{(xv)^k}{k!} \leq |xv|^2 \sum_{k=0}^{\infty} \frac{|xv|^k}{(k+2)!} \leq |xv|^2 \exp(|xv|/2) \leq |xv|^2 \exp\left(\frac{x^2}{4} + \frac{v^2}{4}\right),$$

we may take $h(x) = |x| + 1 + x^2 e^{x^2/4}$, which is certainly integrable against $e^{-\frac{1}{2}x^2}$. \diamond

The tilts that form the basis for lower bounds that do not depend on a particular parameter, but which still provide Fisher-information-like quantities as in Chapter 12.3, also satisfy Definition 13.1.

Example 13.1.8 (Tilted densities and regularity): Let $g : \mathcal{X} \rightarrow \mathbb{R}^d$ be a bounded function, mean zero for a distribution P_0 on \mathcal{X} , and for $t \in \mathbb{R}^d$ define the tilted density

$$p_t(x) := (1 + \langle t, g(x) \rangle) p_0(x)$$

as in definition (12.3.1), where p_0 is the density of P_0 w.r.t. some base measure μ (which can simply be P_0). Then clearly the family $\{P_t\}$ for t in a neighborhood of 0 have densities, and the score $\dot{\ell}_t(x) = \nabla_t \log(1 + tg(x)) = \frac{g(x)}{1 + \langle t, g(x) \rangle}$. Then the ratio $p_t(x)/p_0(x) = 1 + \langle t, g(x) \rangle$, and

$$\left| \frac{p_t(x)}{p_0(x)} - 1 - \langle g(x), t \rangle \right| = 0,$$

satisfying Definition 13.1. \diamond

Then for non-pathological functions f , we typically have the second-order expansion (13.1.7) of the f -divergence in terms of the Fisher information. In this case, we restrict to the class of f -divergences that have global quadratic approximations: we assume there is some $K < \infty$ such that

$$\left| f(1+t) - f(1) - f'(1)t - \frac{f''(1)}{2}t^2 \right| \leq Kt^2 \quad \text{for all } t \geq -1. \quad (13.1.9)$$

The three most common divergences we encounter excepting the variation distance all satisfy this condition.

Example 13.1.9: For the Hellinger distance with $f(t) = \frac{1}{2}(\sqrt{t} - 1)^2$, inequality (13.1.9) holds with $K = \frac{1}{2}$. For the KL-divergence with $f(t) = t \log t - t + 1$, inequality (13.1.9) again holds with $K = \frac{1}{2}$. For the χ^2 divergence with $f(t) = (t - 1)^2$, inequality (13.1.9) holds with $K = 0$. Exercise 13.2 asks you to prove these claims. \diamond

Definition 13.1 and inequality (13.1.9) are enough to guarantee the local information approximation (13.1.7), meaning that for “regular enough” parametric families, all locally quadratic divergences are locally equivalent to the metric induced by the Fisher information matrix.

Lemma 13.1.10. *Let $\{P_\theta\}$ be sufficiently regular for the χ^2 -divergence (Definition 13.1) at θ and $f : \mathbb{R}_+ \rightarrow \mathbb{R}$ be a convex function with $f(1) = 0$ satisfying the second-order Taylor bound (13.1.9). Then for v near 0, the f -divergence satisfies*

$$D_f(P_{\theta+v} \| P_\theta) = \frac{f''(1)}{2} D_{\chi^2}(P_{\theta+v} \| P_\theta) + o(\|v\|^2)$$

and

$$D_{\chi^2}(P_{\theta_0+v} \| P_{\theta_0}) = v^\top J(\theta_0)v + o(\|v\|^2).$$

See Section 13.5.2 for a proof of this result, which as a consequence demonstrates that the χ^2 -divergence is indeed (locally) finite on $\{P_\theta\}$.

JCD Comment: Make this its own subsection to set it off a little bit. Make an example with just the 1-dimensional families. Connect with χ^2 and other divergences here, so that Fisher-informations and other divergences aren't so different. Move Lemma 13.1.10 to here or the next subsection. Add in some results on $(1 + tg)$ families, which will satisfy the regularity conditions and allow “nonparametric” settings, more or less. Also, maybe do the heuristic versions of these, ignoring regularity conditions. Then everything is a bunch cleaner I think.

JCD Comment: Outline for this section: Might be good to actually begin the whole thing with this section.

1. Introduce modulus of continuity (w.r.t. Hellinger), draw a picture suggesting why it should be hard or easy
2. Example with Fisher information-type quantity
3. (Don't do this!) Show that for *testing*, the rate at which we can test really is this modulus whenever we have linear functions and convex classes, because of Le Cam's result on Hellinger affinities.

13.2 Le Cam's convex hull method

JCD Comment: This isn't the starting point any longer. Say that sometimes the method can fail when we just do two point lower bounds, because it might get the wrong rate. Add an exercise about that (e.g., for estimating $\int f'(x)^2 dx$ or something, where motivation might be the type of regularization to use in smoothing the function).

Our starting point is to revisit Le Cam's method from Chapter 9.3, which focused on “two-point” methods to provide a lower bound on estimation error. We can substantially generalize this by instead comparing families of distributions that all induce separations between statistics of one another, and then computing the distance between the convex hulls of the families. This leads to Le Cam's *convex hull method*, which we state abstractly and specialize later to different scenarios of interest. Let \mathcal{P} be a collection of distributions on an underlying space \mathcal{X} , and let $\theta : \mathcal{P} \rightarrow \mathbb{R}^d$ be a parameter of interest. We say that two subsets $\mathcal{P}_0 \subset \mathcal{P}$ and \mathcal{P}_1 are δ -separated in $\|\cdot\|$ if

$$\|\theta(P_0) - \theta(P_1)\| \geq \delta \quad \text{for all } P_0 \in \mathcal{P}_0 \text{ and } P_1 \in \mathcal{P}_1. \quad (13.2.1)$$

We do not require that all of \mathcal{P}_0 be somehow on one side or the other of the collection $\{\theta(P_1) \mid P_1 \in \mathcal{P}_1\}$ of parameters associated with \mathcal{P}_1 , just that they be pairwise separate.

Let $\text{Conv}(\mathcal{P})$ be the collection of mixtures of elements of \mathcal{P} , that is,

$$\text{Conv}(\mathcal{P}) = \left\{ \sum_{i=1}^m \lambda_i P_i \mid m \in \mathbb{N}, \lambda \succeq 0, \langle \lambda, \mathbf{1} \rangle = 1, P_i \in \mathcal{P} \right\}.$$

Defining the minimax risk

$$\mathfrak{M}(\theta(\mathcal{P}), \|\cdot\|) = \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\theta} - \theta(P)\| \right]$$

(note the temporary lack of sample size n), we then have the following generalization of inequality (9.3.3).

Theorem 13.2.1 (Le Cam's Convex Hull Lower Bound). *Let \mathcal{P}_0 and $\mathcal{P}_1 \subset \mathcal{P}$ be δ -separated in $\|\cdot\|$. Then*

$$\mathfrak{M}(\theta(\mathcal{P}), \|\cdot\|) \geq \frac{\delta}{2} \sup \{ [1 - \|\bar{P}_0 - \bar{P}_1\|_{\text{TV}}] \mid \bar{P}_0 \in \text{Conv}(\mathcal{P}_0), \bar{P}_1 \in \text{Conv}(\mathcal{P}_1) \}$$

Proof For any parameter θ , the separation $\|\theta(P_0) - \theta(P_1)\| \geq \delta$ and the triangle inequality guarantees that at least one of $\|\theta - \theta(P_0)\| \geq \delta/2$ or $\|\theta - \theta(P_1)\| \geq \delta/2$ holds for all pairs $P_0 \in \mathcal{P}_0$ and $P_1 \in \mathcal{P}_1$. Let $\bar{P}_0 = \sum_{j=1}^m \alpha_j P_j$ and $\bar{P}_1 = \sum_{j=1}^m \beta_j Q_j$ for $P_j \in \mathcal{P}_0$ and $Q_j \in \mathcal{P}_1$, respectively, where α, β are convex combinations. Then by Markov's inequality,

$$\begin{aligned} \mathfrak{M}(\theta(\mathcal{P}), \|\cdot\|) &\geq \frac{1}{2} \sum_{j=1}^m \alpha_j \mathbb{E}_{P_j} [\|\hat{\theta} - \theta(P_j)\|] + \frac{1}{2} \sum_{j=1}^m \beta_j \mathbb{E}_{Q_j} [\|\hat{\theta} - \theta(Q_j)\|] \\ &\geq \frac{\delta}{2} \left[\sum_{j=1}^m \alpha_j \mathbb{E}_{P_j} [\mathbf{1}\{\|\hat{\theta} - \theta(P_j)\| \geq \delta/2\}] + \sum_{j=1}^m \beta_j \mathbb{E}_{Q_j} [\mathbf{1}\{\|\hat{\theta} - \theta(Q_j)\| \geq \delta/2\}] \right] \\ &\geq \delta \sum_{j=1}^m \left(\alpha_j \mathbb{E}_{P_j} \left[\inf_{P_0 \in \mathcal{P}_0} \mathbf{1}\{\|\hat{\theta} - \theta(P_0)\| \geq \delta/2\} \right] + \beta_j \mathbb{E}_{Q_j} \left[\inf_{P_1 \in \mathcal{P}_1} \mathbf{1}\{\|\hat{\theta} - \theta(P_1)\| \geq \delta/2\} \right] \right) \\ &= \frac{\delta}{2} \left(\mathbb{E}_{\bar{P}_0} \left[\inf_{P_0 \in \mathcal{P}_0} \mathbf{1}\{\|\hat{\theta} - \theta(P_0)\| \geq \delta/2\} \right] + \mathbb{E}_{\bar{P}_1} \left[\inf_{P_1 \in \mathcal{P}_1} \mathbf{1}\{\|\hat{\theta} - \theta(P_1)\| \geq \delta/2\} \right] \right). \end{aligned}$$

Note that if we define $f_v(x) = \inf_{P \in \mathcal{P}_v} \mathbf{1}\{\|\hat{\theta}(x) - \theta(P)\| \geq \delta/2\}$ for $v = 0, 1$, then $f_0 + f_1 \geq 1$. We claim the following lemma, which extends Le Cam's lemma (Proposition 2.3.1) to give

Lemma 13.2.2. *For any two distributions P_0 and P_1 ,*

$$\inf_{f_0+f_1 \geq 1} \mathbb{E}_{P_0}[f_0] + \mathbb{E}_{P_1}[f_1] \geq 1 - \|P_0 - P_1\|_{\text{TV}}.$$

We leave this form of total variation distance as an exercise (see Exercise 2.1). Substituting it into the display above, we find that for any $\bar{P}_v \in \text{Conv}(\mathcal{P}_v)$, we have

$$\mathfrak{M}(\theta(\mathcal{P}), \|\cdot\|) \geq \frac{\delta}{2} [1 - \|\bar{P}_0 - \bar{P}_1\|_{\text{TV}}].$$

Taking a supremum over the \bar{P}_v gives the theorem. \square

13.2.1 The χ^2 -mixture bound

Theorem 13.2.1 provides a powerful tool for developing lower bounds between collections of well-separated distributions. The most typical approach is to take the class \mathcal{P}_0 to consist of a single “base” distribution P_0 , and then let \mathcal{P}_1 vary around P_0 in some prescribed way, so that for an index set \mathcal{V} , we let $\bar{P} = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P_v$. Even so, when we have a sample of size n from one of the distributions, this results in a total variation quantity of the form

$$\|P_0^n - \bar{P}^n\|_{\text{TV}} \quad \text{where} \quad \bar{P}^n = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P_v^n,$$

yielding a mixture of product distributions—something frequently quite challenging to control.

The key technique here is to leverage the inequalities relating divergences from Chapter 2, which allows us to replace the variation distance with something more convenient. In previous chapters, this was the KL-divergence; now, instead, we use a χ^2 -divergence, as it interacts much more nicely with the mixture product structure. Essentially, we replace an expectation over $X \sim P$ with two expectations: one over $X \sim P$ and another over independent samples $V, V' \sim \text{Uniform}(\mathcal{V})$. To obtain the bound, first note that

$$2 \|P_0 - \bar{P}\|_{\text{TV}}^2 \leq D_{\text{kl}}(\bar{P} \| P_0) \leq \log(1 + D_{\chi^2}(\bar{P} \| P_0)) \leq D_{\chi^2}(\bar{P} \| P_0)$$

by Propositions 2.2.8 and 2.2.9.

We then have the following technical lemma.

Lemma 13.2.3. *Let $\bar{P} = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P_v$ and P_v and P_0 have densities p_v, p_0 with respect to some base measure μ on a set \mathcal{X} . Then*

$$D_{\chi^2}(\bar{P} \| P_0) = \frac{1}{|\mathcal{V}|^2} \sum_{v, v' \in \mathcal{V}} \int \frac{p_v(x)p_{v'}(x)}{p_0(x)} d\mu(x) - 1 = \frac{1}{|\mathcal{V}|^2} \sum_{v, v' \in \mathcal{V}} \mathbb{E}_0 \left[\frac{p_v(X)p_{v'}(X)}{p_0^2(X)} \right] - 1,$$

where the expectation is taken with respect to $X \sim P_0$. More generally, let $V \in \mathcal{V}$ be a random variable distributed according to π and conditional on $V = v$, let $X | V = v \sim P_v$. Then for the paired likelihood ratio $l(x | v, v') = \frac{p_v(x)p_{v'}(x)}{p_0^2(x)}$, the marginal distribution \bar{P} of X satisfies

$$D_{\chi^2}(\bar{P} \| P_0) = \mathbb{E}_0 [l(X | V, V')] - 1,$$

where the expectation is taken jointly over $X \sim P_0$ and $V, V' \stackrel{\text{iid}}{\sim} \pi$.

Proof The starting point is to notice that for any two distributions P and Q we have $D_{\chi^2}(P \| Q) = \int (dP/dQ - 1)^2 dQ = \int \frac{dP^2}{dQ} - 2 \int \frac{dP}{dQ} dQ + \int dQ = \int \frac{dP^2}{dQ} - 1$. Then we proceed by recognizing that $(\frac{1}{N} \sum_{i=1}^N x_i)^2 = \frac{1}{N^2} \sum_{i,j} x_i x_j$ for any sequence x_i , and so

$$D_{\chi^2}(\bar{P} \| P_0) + 1 = \int \frac{((1/|\mathcal{V}|) \sum_{v \in \mathcal{V}} dP_v)^2}{dP_0} = \frac{1}{|\mathcal{V}|^2} \sum_{v, v' \in \mathcal{V}} \int \frac{dP_v dP_{v'}}{dP_0}$$

as desired. The second statement has identical proof to the first except that we replace $\frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}}$ with expectations according to π . \square

When we apply Lemma 13.2.3 for product distributions, we can sometimes obtain tensorization-type inequalities, which allows more applications; the next result is an immediate consequence of Lemma 13.2.3.

Lemma 13.2.4. *Let $\bar{P}^n = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P_v^n$, where P_v and P_0 have densities as in Lemma 13.2.3. Then*

$$1 + D_{\chi^2}(\bar{P}^n \| P_0^n) = \frac{1}{|\mathcal{V}|^2} \sum_{v, v' \in \mathcal{V}} \mathbb{E}_0 \left[\frac{p_v(X)p_{v'}(X)}{p_0^2(X)} \right]^n.$$

The applications of these lemmas are many, and going through a few examples will best show how to leverage them. Roughly, our typical approach is the following: we identify \mathcal{V} with $\{\pm 1\}^d$ or some other suitably nice collection of vectors. We then choose distributions P_v and P_0 with densities suitably nice that the ratios p_v/p_0 “act” like exponentials involving inner products of $v \in \mathcal{V}$ with some other quantity; then, because v is uniform in \mathcal{V} in Lemma 13.2.3, we can leverage all the tools we have developed to control moment generating functions and concentration inequalities in Chapter 4 to bound the χ^2 -divergence and then apply Theorem 13.2.1.

Let us give one example of this approach, where we see the technique we use to prove the lemma arises frequently. Let $P_0 = \mathbf{N}(0, \sigma^2 I_d)$ be the standard normal distribution on \mathbb{R}^d , and for $\mathcal{V} = \{-1, 1\}^d$ and some $\delta \geq 0$ to be chosen, let $P_v = \mathbf{N}(\delta v, \sigma^2 I_d)$. Then we have the following lemma, which shows that while $D_{\text{kl}}(P_v \| P_0) = \frac{d\delta^2}{2\sigma^2}$ for each individual P_v , the divergence for the average can be much smaller (even quadratically so in the ratio δ^2/σ^2).

Lemma 13.2.5. *Let P_0 and P_v be Gaussian distributions as above, and define the mixture $\bar{P} = \frac{1}{2^d} \sum_{v \in \{\pm 1\}^d} P_v$. Then*

$$2 \|P_0 - \bar{P}\|_{\text{TV}}^2 \leq \log(1 + D_{\chi^2}(\bar{P} \| P_0)) \leq \frac{d\delta^4}{2\sigma^4}.$$

Proof The first inequality combines Pinsker’s inequality (Proposition 2.2.8) with the bound $D_{\text{kl}}(P \| Q) \leq \log(1 + D_{\chi^2}(P \| Q))$ in Proposition 2.2.9. Now we expand the χ^2 -divergence, yielding

$$1 + D_{\chi^2}(\bar{P} \| P_0) = \mathbb{E} \left[\exp \left(-\frac{1}{2\sigma^2} \|Y - \delta V\|_2^2 - \frac{1}{2\sigma^2} \|Y - \delta V'\|_2^2 + \frac{1}{\sigma^2} \|Y\|_2^2 \right) \right],$$

where the expectation is over $Y \sim \mathbf{N}(0, \sigma^2 I_n)$ and $V, V' \stackrel{\text{iid}}{\sim} \text{Uniform}(\mathcal{V})$. Taking the expectation over Y first, before averaging over the packing elements, allows more careful control. Indeed, expanding the squares and recognizing that $\|v\|_2^2 = d$ for each $v \in \{\pm 1\}^d$, we have

$$\begin{aligned} 1 + D_{\chi^2}(\bar{P} \| P_0) &= \mathbb{E} \left[\exp \left(\frac{\delta}{\sigma^2} \langle Y, V + V' \rangle - \frac{d\delta^2}{\sigma^2} \right) \right] = \mathbb{E} \left[\exp \left(\frac{\delta^2}{2\sigma^2} \|V + V'\|_2^2 - \frac{d\delta^2}{\sigma^2} \right) \right] \\ &= \mathbb{E} \left[\exp \left(\frac{\delta^2}{\sigma^2} \langle V, V' \rangle \right) \right] \\ &\leq \exp \left(\frac{d\delta^4}{2\sigma^4} \right), \end{aligned}$$

where the final key inequality follows because an individual $U \sim \text{Uniform}(\{\pm 1\})$ is 1-sub-Gaussian, and $\langle V, V' \rangle$ is thus d -sub-Gaussian. \square

Using the same techniques, we can provide similar upper bound, coupled with the identity in Lemma 13.2.4, that gives a tensorization-like bound.

Lemma 13.2.6. *Let P_0 and P_v be Gaussian distributions, and define the mixture $\bar{P}^n = \frac{1}{2^d} \sum_{v \in \{\pm 1\}^d} P_v^n$. Then*

$$2 \|P_0^n - \bar{P}^n\|_{\text{TV}}^2 \leq \log(1 + D_{\chi^2}(\bar{P}^n \| P_0^n)) \leq \frac{dn^2\delta^4}{2\sigma^4}.$$

Proof Tracing the proof of Lemma 13.2.5 and using Lemma 13.2.4, we have

$$\begin{aligned} 1 + D_{\chi^2}(\overline{P^n} \| P_0^n) &= \frac{1}{4^d} \sum_{v, v' \in \{\pm 1\}^d} \mathbb{E}_0 \left[\exp \left(\frac{\delta}{\sigma^2} \langle Y, v + v' \rangle - \frac{d\delta^2}{\sigma^2} \right) \right]^n \\ &= \frac{1}{4^d} \sum_{v, v' \in \{\pm 1\}^d} \exp \left(\frac{n\delta^2}{2\sigma^2} \|v + v'\|_2^2 - \frac{nd\delta^2}{\sigma^2} \right). \end{aligned}$$

Letting $V, V' \stackrel{\text{i.i.d.}}{\sim} \text{Uniform}(\{\pm 1\}^d)$ and recognizing that $\|v + v'\|_2^2 = 2d + 2\langle v, v' \rangle$, we have

$$1 + D_{\chi^2}(\overline{P^n} \| P_0^n) = \mathbb{E} \left[\exp \left(\frac{\delta^2 n}{\sigma^2} \langle V, V' \rangle \right) \right] \leq \exp \left(\frac{d\delta^4 n^2}{2\sigma^4} \right)$$

as desired. \square

13.2.2 Estimating the norm of a Gaussian vector

JCD Comment: It would probably be good to connect this to some other literatures and motivate things, e.g.,

1. Signal detection: is there something to discover?
2. Multiple testing: say we have d distinct p-values U_j . Then set $Z_j = \Phi^{-1}(U_j)$. Under the null that $U_j \sim \text{Uniform}[0, 1]$ these are i.i.d. $\mathcal{N}(0, 1)$. Alternatives then deviate from this. Often interesting to consider other alternatives (sparse/dense/etc.)

JCD Comment: Clean this up now, because I moved Lemma 13.2.5 up.

Let us give one example to show how the mixture approach suggested by Lemma 13.2.3 works, along with showing that a more naive approach using the two point method of Chapter 9.3 fails to provide the correct bounds. After this we will further develop the techniques. We motivate the example by considering regression problems, then simplify it to a more stylized and easily workable form. Suppose we wish to estimate the best possible loss achievable in a regression problem,

$$\inf_{\theta} \mathbb{E}[(X^\top \theta - Y)^2].$$

For simplicity, assume that $X \sim \mathcal{N}(0, I_d)$, and that “base” distribution P_0 is simply that $Y \sim \mathcal{N}(0, 1)$, while the alternatives are that $Y = X^\top \theta^* + (1 - \|\theta^*\|_2^2)\varepsilon$, where $\varepsilon \sim \mathcal{N}(0, 1)$ and $\|\theta^*\|_2^2 \leq 1$. In either case we have $Y \sim \mathcal{N}(0, 1)$ marginally, while

$$\inf_{\theta} \mathbb{E}_0[(X^\top \theta - Y)^2] = 1 \quad \text{and} \quad \inf_{\theta} \mathbb{E}_{\theta^*}[(X^\top \theta - Y)^2] = 1 - \|\theta^*\|_2^2,$$

so that estimating the final risk is equivalent to estimating the ℓ_2 -norm $\|\theta^*\|_2^2$.

To make the calculations more palatable, let us assume the simpler *Gaussian sequence model*

$$Y = \theta^* + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_n) \tag{13.2.2}$$

where $\theta^* \in \mathbb{R}^n$ satisfies $\|\theta^*\|_2 \leq r$ for some radius r , and we wish to estimate the statistic

$$T(P) := \|\theta^*\|_2^2.$$

Note that $\mathbb{E}[\|Y\|_2^2] = \|\theta^*\|_2^2 + n\sigma^2$, so that a natural estimator is the debiased quantity

$$T_n := \|Y\|_2^2 - n\sigma^2.$$

Using that $\mathbb{E}[\varepsilon_j^2] = 1$ and $\mathbb{E}[\varepsilon_j^4] = 3$, we then obtain

$$\mathbb{E} \left[\left| T_n - \|\theta^*\|_2^2 \right|^2 \right] = \sum_{j=1}^n \text{Var}((\theta_j^* + \sigma \varepsilon_j)^2) = 2n\sigma^4 + \|\theta^*\|_2^2 \sigma^2 \leq 2n\sigma^4 + r^2 \sigma^2.$$

That is, the family $\mathcal{P}_{\sigma,r}$ defined as Gaussian sequence models (13.2.2) with variance σ^2 and $\|\theta^*\|_2^2 \leq r^2$ satisfies

$$\mathfrak{M}_n(T(\mathcal{P}_{\sigma,r}), |\cdot|) \leq \sqrt{2n\sigma^4 + r^2\sigma^2} \leq \sqrt{2n}\sigma^2 + r\sigma. \quad (13.2.3)$$

We first provide the more naive approach. Suppose that we were to use Le Cam's two-point method to achieve a lower bound in this case. The minimax risk from inequality (9.3.3) shows that (for a numerical constant $c > 0$), if P_0 and P_1 are (respectively) $\mathcal{N}(\theta_0, \sigma^2 I_n)$ and $\mathcal{N}(\theta_1, \sigma^2 I_n)$, then for any choice of θ_0, θ_1 we have

$$\mathfrak{M}_n(T(\mathcal{P}_{\sigma,r}), |\cdot|) \geq \frac{1}{4} \left\{ \left| \|\theta_0\|_2^2 - \|\theta_1\|_2^2 \right| \cdot [1 - \|P_0 - P_1\|_{\text{TV}}] \right\}. \quad (13.2.4)$$

Recalling Pinsker's inequality (Proposition 2.2.8), we have

$$1 - \|P_0 - P_1\|_{\text{TV}} \geq 1 - \frac{1}{\sqrt{2}} \sqrt{D_{\text{kl}}(P_0 \| P_1)} = 1 - \frac{1}{2} \frac{\|\theta_0 - \theta_1\|_2}{\sigma}.$$

So whenever $\|\theta_0 - \theta_1\|_2 \leq \sigma$, we have

$$\mathfrak{M}_n(T(\mathcal{P}_{\sigma,r}), |\cdot|) \geq \frac{1}{8} \left| \|\theta_0\|_2^2 - \|\theta_1\|_2^2 \right|.$$

Take any θ_0 such that $\|\theta_0\|_2 = r$ and $\theta_1 = (1-t)\theta_0$, then choose the largest $t \in [0, 1]$ such that $\|\theta_0 - \theta_1\|_2 = tr \leq \sigma$. The choice $t = \min\{1, \frac{\sigma}{r}\}$ then gives that

$$\|\theta_0\|_2^2 - \|\theta_1\|_2^2 = r^2(1 - (1-t)^2) = r^2(2t - t^2) = 2 \min\{r^2, r\sigma\} - \min\{r^2, \sigma^2\} \geq \min\{r^2, \sigma r\}.$$

In particular, this application of the two-point approach yields

$$\mathfrak{M}_n(T(\mathcal{P}_{\sigma,r}), |\cdot|) \geq \frac{1}{4} \min\{r^2, \sigma r\}. \quad (13.2.5)$$

(A careful inspection of the argument, potentially replacing the application of Pinsker with KL with a Hellinger distance bound, as in Proposition 2.2.8 shows that this is, essentially, the “best possible” bound achievable by the two-point approach.) While this bound *does* capture the second term in the upper bound (13.2.3) whenever $\sigma r \leq r^2$, that is, $r \geq \sigma$, we require more sophisticated techniques to address the scaling with dimension n in the problem.

We therefore turn to using the mixture approach. Let $P_0 = \mathbf{N}(0, \sigma^2 I_n)$, and for $\mathcal{V} = \{\pm 1\}^n$ define $P_v = \mathbf{N}(\delta v, \sigma^2 I_n)$. It is immediate that $T(P_0) = 0$ while $T(P_v) = \delta^2 n$, so we have separation in the values of the statistic. In this case, we apply Theorem 13.2.1 and to obtain

$$\mathfrak{M}_n(T(\mathcal{P}_{\sigma,r}), |\cdot|) \geq \frac{\delta^2 n}{2} \left\{ 1 - \sqrt{\frac{1}{2} \log(1 + D_{\chi^2}(\bar{P} \| P_0))} \right\}$$

for $\bar{P} = \frac{1}{2^n} \sum_{v \in \mathcal{V}} P_v$. Substituting the result of Lemma 13.2.5 into the minimax lower bound, we obtain

$$\mathfrak{M}_n(T(\mathcal{P}_{\sigma,r}), |\cdot|) \geq \frac{\delta^2 n}{2} \left(1 - \sqrt{\frac{n\delta^4}{4\sigma^4}} \right).$$

We choose δ so that the (implied) probability of error in the hypothesis test from which our reduction follows is at least $\frac{1}{2}$, for which it evidently suffices to take $\delta = \frac{\sigma}{n^{1/4}}$. Putting all the pieces together, we achieve the minimax lower bound

$$\mathfrak{M}_n(T(\mathcal{P}_{\sigma,r}), |\cdot|) \geq \frac{\delta^2 n}{4} = \frac{\sigma^2 \sqrt{n}}{4}. \quad (13.2.6)$$

Comparing the result from the upper bound (13.2.3), we see that at least in the regime that the radius r scales at most as $\sigma\sqrt{n}$, the mixture Le Cam method allows us to characterize the minimax risk of estimation of $\|\theta\|_2^2$ in a Gaussian sequence model.

By combining the result (13.2.3) with the more naive two-point lower bound (13.2.5), which is valid in “large radius” regimes, we have actually characterized the minimax risk.

Corollary 13.2.7. *Let $\mathcal{P}_{\sigma,r}$ be the Gaussian sequence model family $\{\mathbf{N}(\theta, \sigma^2 I_n) \mid \|\theta\|_2 \leq r\}$, and $T(\theta) = \|\theta\|_2^2$. Then there is a numerical constant $c > 0$ such that the minimax absolute error satisfies*

$$c(\sigma^2 \sqrt{n} + r\sigma) \leq \mathfrak{M}_n(T(\mathcal{P}_{\sigma,r}), |\cdot|) \leq \sqrt{2n\sigma^4 + r^2\sigma^2}.$$

Proof The only thing to recognize is that $r\sigma \geq \sigma^2 \sqrt{n}$ whenever $r \geq \sigma\sqrt{n}$, in which case $\min\{r^2, \sigma r\} = \sigma r$ in the bound (13.2.5). \square

13.2.3 Lower bounds on estimating integral functionals

When we consider problems such as nonparametric regression or nonparametric density estimation, a frequent concern is the degree of smoothness of the underlying functional, which can then impact the type of regularization one employs—for example, if we know that the underlying function is sufficiently differentiable, cubic smoothing splines estimate nonparametric regression functions by solving

$$\hat{f} = \operatorname{argmin}_f \left\{ \frac{1}{2} \sum_{i=1}^n (Y_i - f(x_i))^2 + \frac{\lambda}{2} \int (f''(x))^2 dx \right\},$$

and the regularization λ is chosen to enforce sufficient smoothness of \hat{f} (see, e.g., [110, Chapter 5.4]).

It is therefore interesting to estimate various functionals based on integration to help understand the gross character of the function being estimated. One can often provide lower bounds on estimation error (of the correct order) for functionals of the form

$$T_k(f) := \int (f^{(k)}(x))^2 dx,$$

the L^2 -norm of the k th derivative of an (unknown) function f , using the convex hull methodology we have so far developed. We will consider estimation of such functionals over classes of functions with bounded higher-order derivatives, defining

$$\mathcal{F}_s := \left\{ f : [0, 1] \rightarrow \mathbb{R} \mid f(0) = 0, f \in \mathcal{C}^s \text{ and } \|f^{(s)}\|_\infty \leq 1 \right\},$$

functions that are s -times continuously differentiable with uniformly bounded s th derivative, where we choose $f(0) = 0$ for normalization. We adopt the observation model of nonparametric regression, as in Example 13.1.3, so that $Y_i = f(X_i) + \varepsilon_i$, where $\varepsilon_i \stackrel{\text{iid}}{\sim} \mathbf{N}(0, 1)$ and $X_i \stackrel{\text{iid}}{\sim} \text{Uniform}([-1, 1])$. Then we can use the convex hull method to obtain lower bounds on estimation.

Proposition 13.2.8. *Fix $k \in \mathbb{N}$. Then for any $s \geq 0$,*

$$\mathfrak{M}_n(T_k(\mathcal{F}_{k+s}), |\cdot|) = \inf_{\hat{T}} \sup_{f \in \mathcal{F}_{k+s}} \mathbb{E}_f \left[|\hat{T}(Y_1^n, X_1^n) - T_k(f)| \right] \gtrsim n^{-\frac{4s}{4(k+s)+1}}.$$

Proof The key is to construct “bump-like” functions that induce separation in the functional T_k , as we have done in the construction of nonparametric lower bounds in Chapter 10.1. Pick any “bump” function $\phi : [0, 1] \rightarrow \mathbb{R}$ with $\phi(0) = \phi(1) = 0$ for which $\|\phi^{(k+s)}\|_\infty \leq 1$ and $\int_0^1 \phi^{(k)}(x)^2 dx > 0$. For example, if we take $h(x) = \exp(-\frac{1}{1-x^2})\mathbf{1}\{|x| < 1\}$, then (up to numerical constant scaling) the bumps $\phi(x) := h(4x - 1) - h(4x - 3)$ satisfy our desiderata, as they are \mathcal{C}^∞ and compactly supported on $[0, 1]$. For a value $m \in \mathbb{N}$ to be chosen, define the rescaled function

$$g(x) := \frac{1}{m^{k+s}} \phi(mx),$$

which evidently satisfies $g^{(k)}(x) = \frac{1}{m^s} \phi^{(k)}(mx)$ and $g^{(k+s)}(x) = \phi^{(k+s)}(mx)$ so $\|g^{(k+s)}\|_\infty \leq 1$. Then for $v \in \{-1, 1\}^m$, define the functions

$$f_v(x) := \sum_{j=1}^m v_j g(mx - (j-1)) = \frac{1}{m^{k+s}} \sum_{j=1}^m v_j \phi(mx - (j-1)),$$

corresponding to bumps of varying directions on each sub-interval $[\frac{j-1}{m}, \frac{j}{m}]$ of $[0, 1]$. We take $f_0(x) = 0$ to be the identically 0 function.

To apply Theorem 13.2.1, the convex hull method, we now follow the usual two steps: we exhibit a separation between the functions f_v and f_0 , and then we bound the divergence between the observations from the null model and the convex hull of the alternatives f_v . The separation is relatively straightforward: clearly we have $T_k(f_0) = 0$, while

$$\begin{aligned} T_k(f_v) &= \int_0^1 f_v^{(k)}(x)^2 dx = m \int_0^{1/m} (g^{(k)}(x))^2 dx = \frac{m}{m^{2s}} \int_0^{1/m} (\phi^{(k)}(mx))^2 dx \\ &= \frac{1}{m^{2s}} \int_0^1 (\phi^{(k)}(u))^2 du \gtrsim \frac{1}{m^{2s}}. \end{aligned} \quad (13.2.7)$$

The divergence bounds for the convex hull are more involved, but follow the approach we have developed: we relate the χ^2 divergence to exponential moments of (independent) random packing vectors, as in Lemma 13.2.5, and then use sub-Gaussianity of random signs (Example 4.1.5) to bound this quantity. Define P_v^n to be the joint distribution of $(X_i, Y_i)_{i=1}^n$ when $Y_i = f_v(X_i) + \varepsilon_i$, and let $\overline{P}^n = \frac{1}{2^m} \sum_{v \in \{\pm 1\}^m} P_v^n$. Let $p_v(y | x)$ denote the density of $Y \sim \mathcal{N}(f_v(x), 1)$. By Lemma 13.2.4, we have

$$\begin{aligned} 1 + D_{\chi^2}(\overline{P}^n \| P_0^n) &= \frac{1}{2^{2m}} \sum_{v, v' \in \{\pm 1\}^m} \mathbb{E}_0 \left[\frac{p_v(Y | X) p_{v'}(X)}{p_0(Y | X)^2} \right]^n \\ &= \frac{1}{2^{2m}} \sum_{v, v' \in \{\pm 1\}^m} \mathbb{E}_0 \left[\exp \left(Y(f_v(X) + f_{v'}(X)) - \frac{1}{2} f_v(X)^2 - \frac{1}{2} f_{v'}(X)^2 \right) \right]^n, \end{aligned}$$

where we have used that $X \sim \text{Uniform}([0, 1])$ and that $Y \sim \mathcal{N}(0, 1)$ under P_0 . Bounding the expectations requires a bit of work. Define the coordinate functions $g_j(x) = g(mx - (j - 1))$ and the vector function $\vec{g}(x) = [g_j(x)]_{j=1}^m$, so that $f_v(x) = \langle v, \vec{g}(x) \rangle$ and $\vec{g}(x)$ has at most one non-zero element, because the supports of the g_j are disjoint. Then

$$\begin{aligned} \mathbb{E}_0 \left[\frac{p_v(Y | X) p_{v'}(X)}{p_0(Y | X)^2} \mid X = x \right] &= \exp \left(\frac{1}{2} \langle v + v', \vec{g}(x) \rangle^2 - \frac{1}{2} \langle v, \vec{g}(x) \rangle^2 - \frac{1}{2} \langle v', \vec{g}(x) \rangle^2 \right) \\ &= \exp(\langle v, \vec{g}(x) \rangle \langle v', \vec{g}(x) \rangle) = \exp(v^\top \text{diag}(\vec{g}(x))^2 v') \end{aligned}$$

because of the disjoint support of the elements of \vec{g} . Now we use that if $|t| \leq 1$, then $e^t \leq 1 + t + t^2$, which implies

$$\begin{aligned} \mathbb{E}_0 \left[\frac{p_v(Y | X) p_{v'}(X)}{p_0(Y | X)^2} \right] &\leq 1 + v^\top \mathbb{E}_0[\text{diag}(\vec{g}(X))^2] v' + \int_0^1 (v^\top \text{diag}(\vec{g}(x))^2 v')^2 dx \\ &= 1 + \int_0^{1/m} g^2(mx) dx \cdot \langle v, v' \rangle + m \int_0^{1/m} g^4(mx) dx \\ &= 1 + \frac{1}{m^{2(k+s)+1}} \int_0^1 \phi^2(u) du \cdot \langle v, v' \rangle + \frac{1}{m^{4(k+s)}} \int_0^1 \phi^4(u) du. \end{aligned}$$

In particular, for numerical constants $c = \int_0^1 \phi^2(u) du$ and $c' = \int_0^1 \phi^4(u) du$, we use that $1 + t \leq e^t$ to obtain that for $V, V' \stackrel{\text{ind}}{\sim} \text{Uniform}(\{\pm 1\}^m)$,

$$\begin{aligned} 1 + D_{\chi^2}(\overline{P}^n \| P_0^n) &\leq \mathbb{E} \left[\exp \left(\frac{nc}{m^{2(k+s)+1}} \langle V, V' \rangle + \frac{nc'}{m^{4(k+s)}} \right) \right] \\ &\leq \exp \left(\frac{c^2 n^2 m}{2m^{4(k+s)+2}} + \frac{nc'}{m^{4(k+s)}} \right). \end{aligned} \tag{13.2.8}$$

In particular, we can choose m scaling as $n^{\frac{2}{4(k+s)+1}}$, which (with an appropriate constant) yields $1 + D_{\chi^2}(\overline{P}^n \| P_0^n) \leq 2$, so that $\|P_0^n - \overline{P}^n\|_{\text{TV}} \leq \log(1 + D_{\chi^2}(\overline{P}^n \| P_0^n)) \leq \log 2 < 1$. Substituting this choice of m into the separation bound (13.2.7) gives the result. \square

Proposition 13.2.8 gives a few consequences. First, if we wish to estimate the integral of the k th derivative of a function f , if all we know is that f has bounded and continuous k th derivative

then the minimax risk is constant—estimation is impossible. Typical choices are that $k = 1$ and the number of additional degrees of smoothness $s = 1$ in Proposition 13.2.8, which implies a lower bound on the minimax risk of

$$\mathfrak{M}(T_1(\mathcal{F}_2), |\cdot|) \gtrsim n^{-\frac{1}{2+1/4}} = n^{-\frac{4}{9}}.$$

More generally, if the number of additional degrees of smoothness $s \leq k$, then standard parametric rates of $n^{-1/2}$ are unachievable.

JCD Comment: Maybe add two or three exercises around these ideas:

1. Failure of the 2-point bound to achieve Proposition 13.2.8
2. The density versions of these arguments
3. The curse of dimensionality appearing even in functional estimation

Also add a figure probably.

13.3 Minimax hypothesis testing

In the general hypothesis testing problem, we have a family of potential distributions \mathcal{P} , and we are given a sample $X \sim P$ for some $P \in \mathcal{P}$. Then we wish to distinguish between two disjoint hypotheses H_0 and H_1 :

$$\begin{aligned} H_0 : & P \in \mathcal{P}_0 \\ H_1 : & P \in \mathcal{P}_1, \end{aligned} \tag{13.3.1}$$

where the collections $\mathcal{P}_0 \subset \mathcal{P}$ and $\mathcal{P}_1 \subset \mathcal{P}$ are disjoint. Then for a given test statistic $\Psi : \mathcal{X} \rightarrow \{0, 1\}$, we define the *risk* of the test to be

$$R(\Psi | \mathcal{P}_0, \mathcal{P}_1) := \sup_{P \in \mathcal{P}_0} P(\Psi \neq 0) + \sup_{P \in \mathcal{P}_1} P(\Psi \neq 1),$$

that is, the sum of the worst-case probabilities that the test is incorrect. (We also use the notation $R(\Psi | H_0, H_1)$ to denote the same quantity.) In the scenarios we consider, we will assume a metric ρ on the family of distributions \mathcal{P} , and instead of the general hypothesis test (13.3.1), we will consider testing whether $P \in \mathcal{P}_0$ or $\rho(P, P_0) \geq \epsilon$ for all $P_0 \in \mathcal{P}_0$, giving the variant

$$\begin{aligned} H_0 : & P \in \mathcal{P}_0 \\ H_1 : & P \in \mathcal{P}_1(\epsilon) := \{P \in \mathcal{P} \text{ s.t. } \rho(P, P_0) \geq \epsilon \text{ all } P_0 \in \mathcal{P}_0\} \end{aligned} \tag{13.3.2}$$

In this case, we can define the *risk at distance ϵ* for a sample of size n by

$$R_n(\Psi, \epsilon) := \sup_{P \in \mathcal{P}_0} P(\Psi(X_1^n) \neq 0) + \sup_{P \in \mathcal{P}_1(\epsilon)} P(\Psi(X_1^n) \neq 1), \tag{13.3.3}$$

leaving \mathcal{P}_0 and \mathcal{P} implicit in the definition, and where we let $X_1^n \stackrel{\text{iid}}{\sim} P$. From this, we can define the minimax test risk

$$\inf_{\Psi} R_n(\Psi, \epsilon).$$

We then ask for the particular thresholds ϵ at which the minimax test risk becomes small or large. Thus, while the coming definition allows some ambiguity, we say that a sequence ϵ_n is a

minimax threshold or *critical testing radius* for the testing problem (13.3.2) if there exist numerical constants $0 < c \leq C < \infty$ such that

$$\inf_{\Psi} R_n(\Psi, C\epsilon_n) \leq \frac{1}{3} \quad \text{and} \quad \inf_{\Psi} R_n(\Psi, c\epsilon_n) \geq \frac{2}{3}. \quad (13.3.4)$$

The constants $\frac{1}{3}$ and $\frac{2}{3}$ are unimportant, the point being that for separation at most $c\epsilon_n$, *no* hypothesis test can test whether the distribution P satisfies $P \in \mathcal{P}_0$ or $\inf_{P_0 \in \mathcal{P}_0} \rho(P, P_0) \geq c\epsilon_n$ with reasonable accuracy. But it *is* possible to test whether $P \in \mathcal{P}_0$ or $\inf_{P_0 \in \mathcal{P}_0} \rho(P, P_0) \geq C\epsilon_n$ with reasonable accuracy. Moreover, we can make the probability of error exponentially small by increasing the sample size by a constant factor, as Exercise 13.4 explores. In some cases, and we give one extended example in Section 13.3.4, one can establish a stronger result than the critical radius (13.3.4), instead establishing a phase transition. In this case, we say that a sequence ϵ_n is the *phase transition threshold* if for any $c < 1 < C$,

$$\limsup_n \inf_{\Psi} R_n(\Psi, C\epsilon_n) = 0 \quad \text{and} \quad \liminf_n \inf_{\Psi} R_n(\Psi, c\epsilon_n) = 1. \quad (13.3.5)$$

Conveniently, the minimax test risk has a precise divergence-based form, to which we can apply the techniques comparing different divergences we have developed. In particular, we have the following analogue of Le Cam's convex hull lower bound in Theorem 13.2.1, which provides the same fundamental quantity (the variation distance between convex hulls of \mathcal{P}_0 and \mathcal{P}_1) for lower bounds, except that it applies for testing.

Proposition 13.3.1 (Convex hull lower bounds in testing). *For any classes \mathcal{P}_0 and \mathcal{P}_1 , the minimax test risk satisfies*

$$\inf_{\Psi} R(\Psi \mid \mathcal{P}_0, \mathcal{P}_1) \geq 1 - \sup \left\{ \|\bar{P}_0 - \bar{P}_1\|_{\text{TV}} \mid \bar{P}_0 \in \text{Conv}(\mathcal{P}_0), \bar{P}_1 \in \text{Conv}(\mathcal{P}_1) \right\}.$$

Proof Let $\bar{P}_0 \in \text{Conv}(\mathcal{P}_0)$ and $\bar{P}_1 \in \text{Conv}(\mathcal{P}_1)$. Then for any test Ψ ,

$$R(\Psi \mid \mathcal{P}_0, \mathcal{P}_1) \geq \bar{P}_0(\Psi \neq 0) + \bar{P}_1(\Psi \neq 1)$$

because suprema are always at least as large as averages. Now note that the set $A = \{x \mid \Psi(x) = 0\}$ satisfies

$$\bar{P}_0(\Psi \neq 0) + \bar{P}_1(\Psi \neq 1) = \bar{P}_0(A^c) + \bar{P}_1(A) = 1 - (\bar{P}_0(A) - \bar{P}_1(A)),$$

and take an infimum over regions A . □

In fact, equality typically holds in Proposition 13.3.1, but this requires the application of (infinite dimensional) convex duality, which is beyond our scope here.

13.3.1 Detecting a difference in populations

With the generic worst-case hypothesis testing setup in place, we can give a general recipe for developing tests. We specialize this recipe in the next few sections to different problems, including signal detection in a Gaussian model, two-sample tests in multinomials, and goodness of fit testing. The basic approach in all of these problems is frequently the following: to demonstrate achievability and testability, we develop an estimator T_n of the distance $\rho(P_0, P_1)$, or some other function of the distance, where T_n has reasonable properties. We then develop a test Ψ by thresholding this

estimator. For the converse results that no test can distinguish the families \mathcal{P}_0 and \mathcal{P}_1 at a particular distance, we use the mixture χ^2 approaches we have outlined.

Let us give the general recipe first. Suppose that we have a statistic T designed to separate the classes \mathcal{P}_0 and \mathcal{P}_1 . Such a statistic should assign large values for samples $X \sim P_1$ for $P_1 \in \mathcal{P}_1$ and small values for samples $X \sim P_0$. A more quantitative version of this, where the separation $\mathbb{E}_1[T] - \mathbb{E}_0[T]$ is commensurate with the variance of T , is sufficient to test between \mathcal{P}_0 and \mathcal{P}_1 with high accuracy. To that end, we say that the statistic T *robustly C -separates* \mathcal{P}_0 and \mathcal{P}_1 if

$$\mathbb{E}_{P_1}[T] - \sup_{P_0 \in \mathcal{P}_0} \mathbb{E}_{P_0}[T] \geq C \left(\sup_{P_0 \in \mathcal{P}_0} \sqrt{\text{Var}_{P_0}(T)} + \sqrt{\text{Var}_{P_1}(T)} \right). \quad (13.3.6)$$

for each $P_1 \in \mathcal{P}_1$. Typically, we choose statistics T so that $\mathbb{E}_{P_0}[T] = 0$ for each P_0 in the null \mathcal{P}_0 (though this is not always possible). The next proposition shows how to define a test that leverages this to achieve small worst-case test error.

Proposition 13.3.2. *Let the statistic $T : \mathcal{X} \rightarrow \mathbb{R}$ robustly C -separate \mathcal{P}_0 from \mathcal{P}_1 . Then for the threshold $\tau = \sup_{P_0 \in \mathcal{P}_0} \mathbb{E}_{P_0}[T] + \sup_{P_0 \in \mathcal{P}_0} \sqrt{\text{Var}_{P_0}(T)}$, the test*

$$\Psi(X) := \mathbf{1}\{T \geq \tau\}$$

satisfies

$$R(\Psi \mid \{P_0\}, \mathcal{P}_1) \leq \frac{2}{C^2}.$$

Proof Without loss of generality we assume $\sup_{P_0 \in \mathcal{P}_0} \mathbb{E}_{P_0}[T] = 0$, as the test is invariant to shifts, so that $\tau = \sup_{P_0 \in \mathcal{P}_0} \sqrt{\text{Var}_{P_0}(T)}$. We can also assume that $C \geq 1$, as otherwise the proposition is vacuous. We control the test error in each case. Under any null P_0 , we have

$$P_0(\Psi \neq 0) = P_0(T \geq \tau) \leq \frac{\text{Var}_0(T)}{C^2 \tau^2} = \frac{1}{C^2}.$$

For the alternatives under $P_1 \in \mathcal{P}_1$, we have

$$P_1(\Psi \neq 1) = P_1(T \leq \tau) = P_1(T - \mathbb{E}_1[T] \leq \tau - \mathbb{E}_1[T]) \leq \frac{\text{Var}_1(T)}{[\mathbb{E}_1[T] - \tau]_+^2}.$$

But of course,

$$\mathbb{E}_1[T] - \tau = \mathbb{E}_1[T] - \sup_{P_0} \mathbb{E}_{P_0}[T] - \sup_{P_0} \sqrt{\text{Var}_{P_0}(T)} \geq C \sqrt{\text{Var}_1(T)} + (C - 1) \sup_{P_0} \sqrt{\text{Var}_{P_0}(T)}$$

by the robust C -separation. As we have assumed w.l.o.g. that $C \geq 1$, this yields

$$P_1(\Psi \neq 1) \leq \frac{\text{Var}_1(T)}{C^2 \text{Var}_1(T)} = \frac{1}{C^2}$$

as desired. □

13.3.2 Signal detection and testing a Gaussian mean

A common problem in statistics, communication, and information theory is the *signal detection* problem, where we observe $X \sim P$ from an unknown distribution P , and wish to detect if there is some “signal” present in P . To study such a problem, we typically formulate a null model—indicating absence of signal—and a set of alternatives for which there is *some* signal, though we only care to test its existence. The existence of a signal can then justify further investigation or data collection to actually estimate the signal.

Let us give a few variants of this problem, for which a substantial literature exists.

Example 13.3.3 (Dense Gaussian signal detection): We consider testing the null H_0 and alternative H_1 given by

$$\begin{aligned} H_0 : P &= P_0 = \mathbf{N}(0, I_d) \\ H_1 : P &\in \mathcal{P}_1(r) := \{\mathbf{N}(\theta, I_d) \mid \|\theta\|_2 \geq r\}. \end{aligned} \quad (13.3.7)$$

That is, we are interested in whether $X \sim P$ has a mean θ separated by at least r from the all-zeros vector. The problem is to find the critical radius r at which testing between $\mathcal{P}_0 = \{P_0\}$ and \mathcal{P}_1 becomes feasible (or infeasible). \diamond

Example 13.3.4 (A global null in multiple hypothesis testing): Consider the problem of testing d distinct null hypotheses $H_{0,j}$, $j = 1, \dots, d$, where for each we have a p -value Y_j and reject $H_{0,j}$ if $Y_{0,j} \leq \tau$ for a threshold τ . (Recall that a p value is a random variable Y that is *sub-uniform*, meaning that $P(Y \leq u) \leq P(U \leq u)$ for $U \sim \text{Uniform}[0, 1]$, so we are less likely to reject at threshold τ than a uniform would be.) If we assume the Y_j are exact p -values, that is, $P(Y_j \leq u) = u$ for $u \in [0, 1]$, then testing the global independent null

$$H_0 := \bigcap_{j=1}^d H_{0,j} = \text{each } Y_j \stackrel{\text{iid}}{\sim} \text{Uniform}[0, 1]$$

is equivalent to Gaussian signal detection. Indeed, let $Z_j = \Phi^{-1}(Y_j)$, where Φ denotes the standard Gaussian cumulative distribution. Then under the global null H_0 , we have

$$Z \sim \mathbf{N}(0, I_d).$$

The question of which alternative class \mathcal{P}_1 to consider is then frequently a matter of applications. For example, we might be curious about alternatives for which a few nulls $H_{0,j}$ are false, that is, *sparse* alternatives. Example 13.3.3 corresponds to something like dense alternatives. \diamond

With these as motivation, let us consider Example 13.3.3 more carefully, in effort to find the critical radius r at which minimax testing becomes feasible (or infeasible). While our standard techniques for estimation tell us that the minimax rate for estimating θ in a normal location family $\mathcal{P} = \{\mathbf{N}(\theta, \sigma^2 I_d)\}_{\theta \in \mathbb{R}^d}$ (say, in mean squared error) necessarily scale as

$$\mathfrak{M}_n(\theta(\mathcal{P}), \|\cdot\|_2^2) = \frac{d\sigma^2}{n},$$

we can *test* whether the mean of a Gaussian is zero at a smaller dimensionality—effectively, while $\mathbb{E}[\|\hat{\theta} - \theta\|_2^2] \rightarrow 0$ as $n \rightarrow \infty$ if and only if $d/n \rightarrow 0$, in the testing case, we can save a dimension-dependent factor \sqrt{d} . In particular, the next two examples—one addressing achievability and one

the fundamental limit—show that in the dense Gaussian signal detection problem of Example 13.3.3, the critical test radius (13.3.4) at which testing is feasible or infeasible scales as

$$r_n := \frac{d^{1/4}}{\sqrt{n}}.$$

We can achieve (asymptotically) accurate testing in the dense signal detection problem (13.3.7) if and only if $\sqrt{d}/n \rightarrow 0$ as $n \rightarrow \infty$.

We first demonstrate achievability in Example 13.3.3, leveraging Proposition 13.3.2.

Example 13.3.5 (Achievability in Gaussian mean testing): We wish to test the alternatives (13.3.7). We use the approach of Proposition 13.3.2: find an estimator of $\|\theta\|_2^2$, and then threshold it for our test. The discussion preceding Corollary 13.2.7 (specifically equation (13.2.3)) shows that given a sample of size n , the estimator $T_n = \|\bar{X}_n\|_2^2 - d/n$ is unbiased for $\|\theta\|_2^2$ and satisfies

$$\mathbb{E}_\theta [(T_n - \|\theta\|_2^2)^2] = \text{Var}_\theta(T_n) \leq \frac{2d}{n^2} + \frac{\|\theta\|_2^2}{n}. \quad (13.3.8)$$

Note that $\mathbb{E}_0[T_n] = 0$, and so because

$$\mathbb{E}_\theta[T_n] - \mathbb{E}_0[T_n] = \|\theta\|_2^2,$$

the statistic T_n robustly 2-separates P_0 from $\mathcal{P}_1(r)$ (recall definition (13.3.6)) whenever

$$\frac{1}{2} \|\theta\|_2^2 \geq \left(\frac{\sqrt{2d}}{n} + \sqrt{\frac{2d}{n^2} + \frac{1}{n} \|\theta\|_2^2} \right)$$

for all θ with $\|\theta\|_2 \geq r$. Immediately we see that if we take radius $r^2 = C\sqrt{d}/n$ for some $C > 0$, then this separation occurs if $C\sqrt{d} \geq 2(\sqrt{2d} + \sqrt{2d + C\sqrt{d}})$, which of course happens for large constant C . Applying Proposition 13.3.2, we thus see that the test $\Psi(X_1^n) = \mathbf{1} \{T_n \geq \sqrt{2d/n^2}\}$ satisfies

$$R_n(\Psi, Cr_n) \leq \frac{1}{3} \quad \text{for } r_n = \frac{d^{1/4}}{\sqrt{n}},$$

which gives the achievability required for the critical test radius (13.3.4). \diamond

Example 13.3.5 shows that at the critical radius $r_n = \frac{d^{1/4}}{\sqrt{n}}$, it is possible (in a worst-case sense) to test between the null $H_0 : \mathbf{N}(0, I_d)$ and alternatives $H_1 : \mathbf{N}(\theta, I_d)$ for $\|\theta\|_2 \geq Cr_n$, where C is a numerical constant. We can also provide the converse.

Example 13.3.6 (Lower bounds in Gaussian mean testing): Let $\mathcal{P}_1(r) = \{\mathbf{N}(\theta, I_d) \mid \|\theta\|_2 \geq r\}$ be a collection of Gaussians with means r away from the origin in ℓ_2 -norm. We seek the critical radius r below which it is impossible to distinguish between $P_0 = \mathbf{N}(0, I_d)$ and $P_1 \in \mathcal{P}_1(r)$ given an i.i.d. sample X_1^n . Lemma 13.2.6 and Proposition 13.3.1 combine (set $\delta^2 = \frac{r^2}{d}$ in Lemma 13.2.6) to give

$$\inf_{\Psi} R_n(\Psi \mid \mathcal{P}_0, \mathcal{P}_1(r)) \geq 1 - \sqrt{\frac{n^2 r^4}{4d}}.$$

In particular, the threshold $r^2 = \sqrt{d}/n$ means that there is necessarily constant test error probability $R_n \geq \frac{1}{2}$. Combining the estimation guarantee with this lower bound shows that the critical radius (13.3.4) for testing $H_0 : \mathbf{N}(0, I_d)$ against the family of alternatives $H_1 : \mathbf{N}(\theta, I_d)$ with $\|\theta\|_2^2 \geq r^2$ is precisely $r^2 = \sqrt{d}/n$. \diamond

13.3.3 Goodness of fit and two-sample tests for multinomials

The basic question in goodness of fit testing—called property testing in the theoretical computer science literature—is the following. Given a sample $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P$, we wish to test whether $P = P_0$ for a prescribed base distribution P_0 or P is far from P_0 . The related two-sample testing problem generalizes this, where we assume samples $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P$ and $Y_1, \dots, Y_m \stackrel{\text{iid}}{\sim} Q$, and wish to test whether $P = Q$. Each of these falls into the class of hypothesis tests (13.3.2), where the choice of the metric ρ can change the character of upper and lower bounds somewhat dramatically. General methods for developing goodness of fit and two-sample tests typically take the broad approach in Section 13.3.1, defining a statistic T that separates the distribution P_0 (or the joint that X_i and Y_j have the same distribution) from the alternatives about which we are curious, then thresholding that statistic.

It turns out that even in what might appear to be a particularly simple case—that of multinomial distributions, where we identify the distribution P with a probability mass function (p.m.f.) $p \in \Delta_d$ —a surprising amount of complexity arises. We thus work through two examples on testing distance between discrete distributions by considering two metrics on the probability mass functions: the ℓ_2 -metric and the total variation distance (or ℓ_1 metric). Then $\rho(p, q) = \|p - q\|$ for $\|\cdot\| = \|\cdot\|_2$ or $\|\cdot\| = \|\cdot\|_1$. In the uniformity testing case, we let $p_0 = \frac{1}{d}\mathbf{1}$ be the uniform distribution on $[d]$, and we seek the critical threshold ϵ at which testing

$$\|p - p_0\| = 0 \quad \text{versus} \quad \|p - p_0\| \geq \epsilon$$

from n i.i.d. observations $X_i \stackrel{\text{iid}}{\sim} p$ becomes feasible or infeasible.

It is simpler (for analyzing procedures) to consider a slight variant of this problem, which uses the *Poissonization* trick. To motivate the idea, identify the observations X_i with the basis vectors (so that observing item $j \in \{1, \dots, d\}$ corresponds to $X_i = e_j$). Then that the sample mean $\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i$ is unbiased, but its coordinates exhibit dependence in that $\langle \mathbf{1}, \hat{p} \rangle = 1$ —an annoyance for analyses. Thus, we consider an alternative approach, where we assume a two-stage sampling procedure: we first draw $N \sim \text{Poi}(n)$, and then conditional on $N = m$, draw $X_i \stackrel{\text{iid}}{\sim} p$, $i = 1, \dots, m$. As $\mathbb{E}[N] = n$ and N concentrates around its mean, this is nearly equivalent to simply observing $X_i \stackrel{\text{iid}}{\sim} p$ for $i = 1, \dots, n$, and a standard probabilistic calculation shows that the distribution of $\{X_i\}_{i=1}^N$ conditional on $N = m$ is identical to the distribution of $X_i \stackrel{\text{iid}}{\sim} p$, $i = 1, \dots, m$.

Even more, the minimax risk for estimation in this Poissonized sampling scheme is similar to that for estimation in the original multinomial setting. Indeed, suppose that we wish to estimate an abstract statistic $T(p)$ of $p \in \Delta_d$, and assume for simplicity that $T(p) \in [-r, r]$ for some fixed r . Define the minimax and Poissonized minimax risks

$$\mathfrak{M}_n := \inf_{T_n} \sup_{p \in \Delta_d} \mathbb{E}_p [(T_n(X_1^n) - T(p))^2]$$

and

$$\mathfrak{M}_{\text{Poi}(n)} := \inf_{\{T_m\}} \sup_{p \in \Delta_d} \mathbb{E}_p [(T_N(X_1^N) - T(p))^2],$$

where the latter expectation is taken over the sample size $N \sim \text{Poi}(n)$, and $\{T_m\}$ denotes a sequence of estimators (defined for all sample sizes m). We have the following proposition, which shows that if we can provide procedures that work in the poissonized (independent sampling) setting, then the standard multinomial sampling setting is similarly easy (or challenging).

Proposition 13.3.7. *There exist numerical constants $0 < c, C < \infty$ such that*

$$\mathfrak{M}_{\text{Poi}(2n)} - Cr^2 \exp(-cn) \leq \mathfrak{M}_n \leq 2 \cdot \mathfrak{M}_{\text{Poi}(n/2)}. \quad (13.3.9)$$

For a proof, see Exercises 13.12 and 13.13.

Let us leverage these ideas to construct an estimator for the ℓ_2 -distance between two multinomial distributions. In this case, suppose we have $X_i \stackrel{\text{iid}}{\sim} p$ and $Y_i \stackrel{\text{iid}}{\sim} q$, where $p, q \in \Delta_d$, both for $i = 1, \dots, N$ and $N \sim \text{Poi}(n)$, and we define

$$\hat{p} = \frac{1}{n} \sum_{i=1}^N X_i, \quad \hat{q} = \frac{1}{n} \sum_{i=1}^N Y_i. \quad (13.3.10)$$

This is equivalent to sampling $n\hat{p}_j \stackrel{\text{ind}}{\sim} \text{Poi}(np_j)$ and $n\hat{q}_j \stackrel{\text{ind}}{\sim} \text{Poi}(nq_j)$, $j = 1, \dots, d$, and so we use the quantities (13.3.10) to define an estimator we can threshold using Proposition 13.3.1. We work through this in the next (somewhat complicated) example.

Example 13.3.8 (Estimating the ℓ_2 -distance between multinomials): For the estimators (13.3.10), define the quantity

$$Z_j := (n\hat{p}_j - n\hat{q}_j)^2 - n\hat{p}_j - n\hat{q}_j.$$

Recalling that if $W \sim \text{Poi}(\lambda)$ then $\mathbb{E}[W] = \text{Var}(W) = \lambda$, we have $\mathbb{E}[n\hat{p}_j] = p_j$ and $\text{Var}(n\hat{p}_j) = np_j$, so

$$\begin{aligned} \mathbb{E}[Z_j] &= \mathbb{E}[(n\hat{p}_j)^2] + \mathbb{E}[(n\hat{q}_j)^2] - 2n^2 p_j q_j - np_j - nq_j \\ &= \text{Var}(n\hat{p}_j) + \text{Var}(n\hat{q}_j) + (np_j)^2 + (nq_j)^2 - 2n^2 p_j q_j - np_j - nq_j = n^2 \|p - q\|_2^2. \end{aligned}$$

In particular, the statistic

$$T_n := \frac{1}{n^2} \langle \mathbf{1}, Z \rangle$$

satisfies $\mathbb{E}[T_n] = \|p - q\|_2^2$.

To be able to test whether p and q are identical using Proposition 13.3.2, we must compute the variance of $\langle \mathbf{1}, Z \rangle$, which—conveniently, by the independence our Poisson sampling gives—is $\sum_{j=1}^d \text{Var}(Z_j)$. Leveraging that for a Poisson $W \sim \text{Poi}(\lambda)$ we have (by tedious calculation) that

$$\mathbb{E}[W] = \lambda, \quad \mathbb{E}[W^2] = \lambda(1 + \lambda), \quad \mathbb{E}[W^3] = \lambda + 3\lambda^2 + \lambda^3, \quad \mathbb{E}[W^4] = \lambda + 7\lambda^2 + 6\lambda^3 + \lambda^4,$$

we obtain (see Exercise 13.16)

$$\text{Var}(Z_j) = 4n^3(p_j - q_j)^2(p_j + q_j) + 2(p_j + q_j)^2 n^2 \quad (13.3.11)$$

and

$$\text{Var}(\langle \mathbf{1}, Z \rangle) \leq 4n^3 \|p - q\|_4^2 \|p + q\|_2 + 2n^2 \|p + q\|_2^2.$$

Under the (non-point) null $H_0 : p = q$, $\text{Var}(\langle \mathbf{1}, Z \rangle) = 2n^2 \|p + q\|_2^2 \leq 8n^2$, as $\sup_{p,q} \|p + q\|_2 = 2$. Proposition 13.3.2 thus shows that if

$$\|p - q\|_2^2 \geq C \left(\sqrt{\frac{8}{n^2}} + \sqrt{\frac{16 \|p - q\|_4^2}{n} + \frac{8}{n^2}} \right), \quad (13.3.12)$$

then the test

$$\Psi := \mathbf{1} \left\{ T_n \geq \sqrt{8/n} \right\}$$

satisfies $P_0(\Psi \neq 0) + P_1(\Psi \neq 1) \leq \frac{2}{\sqrt{d}}$, where P_0 is any distribution with $p = q$ and P_1 is any distribution with $\|p - q\|_2$ satisfying the separation (13.3.12). As $\|p - q\|_2 \geq \|p - q\|_4$, inequality (13.3.12) a necessary and sufficient condition for inequality (13.3.12) to hold is that $\|p - q\|_2 \gtrsim 1/\sqrt{n}$. \diamond

Summarizing, we see that if we wish to test whether two multinomials are identical or separated in ℓ_2 , the critical threshold for the hypothesis test

$$\begin{aligned} H_0 : & \quad p = q \\ H_1 : & \quad \|p - q\|_2 \geq \delta \end{aligned} \tag{13.3.13}$$

satisfies $\delta \leq \frac{1}{\sqrt{n}}$: we can test between H_0 and H_1 at separations that are essentially “independent” of the dimension or number of categories d . This is in fact sharp, as a relatively straightforward argument with Le Cam’s two-point lemma demonstrates (see Exercise 13.18). However, if we change the norm $\|\cdot\|_2$ into the ℓ_1 -norm $\|\cdot\|_1$, the story changes significantly.

Let us change the hypothesis test (13.3.13) to simpler looking—in that we only test goodness of fit— ℓ_1 -based variant. Identifying distributions P on $\{1, \dots, d\}$ with their p.m.f.s $p \in \Delta_d$, let P_0 be the uniform distribution on $\{1, \dots, d\}$, with p.m.f. $p_0 = \frac{1}{d}\mathbf{1}$. Then we consider the testing problem

$$\begin{aligned} H_0 : & \quad p = p_0 \\ H_1 : & \quad \|p - p_0\|_1 \geq \delta, \end{aligned} \tag{13.3.14}$$

which tests the ℓ_1 -distance to uniformity. In this case, developing a test that distinguishes these hypotheses at the optimal rate is quite sophisticated, though we outline an approach to it in the exercises. To develop the correct order of lower bound—that is, a threshold δ for which no test can reliably distinguish H_0 from H_1 —is possible via the mixture of χ^2 -distributions approach we have developed in Lemma 13.2.3.

Proposition 13.3.9 (A lower bound for testing ℓ_1 -separated multinomials). *In the testing problem (13.3.14),*

$$\inf_{\Psi} R_n(\Psi \mid H_0, H_1) \geq 1 - \frac{1}{\sqrt{2}}$$

whenever $\delta \leq \frac{d^{1/4}}{\sqrt{n}}$.

Proof We construct a particular packing of the probability simplex $\Delta_d \in \mathbb{R}_+^d$ that guarantees that the divergence between elements of H_0 and H_1 in the test (13.3.14) is small. For simplicity, we assume d is even, as it changes nothing. For the base distribution P_0 take p.m.f. $p_0 = \frac{1}{d}\mathbf{1}$ as required by the problem (13.3.14). To construct the alternatives, let $\mathcal{V} \subset \{\pm 1\}^d$ be the collection of $2^{d/2}$ vectors of the form $v = (v', -v')$, where $v' \in \{\pm 1\}^{d/2}$, so that $\langle \mathbf{1}, v \rangle = 0$ for each $v \in \mathcal{V}$. Then for $\delta \geq 0$ to be chosen, define the p.m.f.s $p_v = \frac{1+\delta v}{d}$. Identify samples $X \in \{e_1, \dots, e_d\}$. Then for any $x \in \{e_j\}$, we have $P_v(X = x) = \frac{1}{d}(1 + \delta \langle v, x \rangle)$, and so for any pair v, v' we have

$$\frac{P_v(X = x)P_{v'}(X = x)}{P_0(X = x)^2} = (1 + \delta \langle v, x \rangle)(1 + \delta \langle v', x \rangle).$$

From this key equality, we see that if $V, V' \stackrel{\text{iid}}{\sim} \text{Uniform}(\mathcal{V})$, then for $\bar{P} = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P_v$ we have

$$\begin{aligned} 1 + D_{\chi^2}(\bar{P} \| P_0) &= \mathbb{E}_0 \left[\prod_{i=1}^n (1 + \delta \langle V, X_i \rangle) (1 + \delta \langle V', X_i \rangle) \right] \\ &= \mathbb{E} \left[\mathbb{E}_0[(1 + \delta \langle V, X \rangle)(1 + \delta \langle V', X \rangle) \mid V, V']^n \right] \\ &= \mathbb{E} \left[\left(1 + \frac{\delta^2}{d} \langle V, V' \rangle \right)^n \right], \end{aligned}$$

where the final equality follows because $\mathbb{E}_0[\langle v, X \rangle] = \frac{1}{d} \langle v, \mathbf{1} \rangle = 0$ for each $v \in \mathcal{V}$. Now we use that $1 + t \leq e^t$ for all t to obtain

$$1 + D_{\chi^2}(\bar{P} \| P_0) \leq \mathbb{E} \left[\exp \left(\frac{n\delta^2}{d} \langle V, V' \rangle \right) \right] = \mathbb{E} \left[\exp \left(\frac{2n\delta^2}{d} \sum_{j=1}^{d/2} U_j \right) \right]$$

for $U_j \stackrel{\text{iid}}{\sim} \text{Uniform}\{\pm 1\}$. But of course these U_j are 1-sub-Gaussian, so

$$1 + D_{\chi^2}(\bar{P} \| P_0) \leq \exp \left(\frac{n^2 \delta^4}{d} \right).$$

Now use Pinsker's inequalities (Propositions 2.2.8 and 2.2.9), which gives $2 \|P_0 - \bar{P}\|_{\text{TV}}^2 \leq \frac{n^2}{\delta^4} d$. Choose $\delta^4 = \frac{d}{n^2}$. \square

13.3.4 Detecting sparse signals and phase transitions

Interesting phenomena arise when we consider signal detection problems, as in Section 13.3.2, but the underlying signal is sparse. For example, in an astronomical survey, where we search for light sources, we may be interested in regions of space where there are more than typical number of astronomical objects. Then the signal is quite sparse—much of space is empty—but we still wish to make discoveries. We can distill much of the complexity here into an example motivating much of our development.

Example 13.3.10 (Sparse Gaussian signal detection): In the sparse Gaussian signal detection problem, we observe n random variables $Y_i \sim \mathbf{N}(\mu_i, 1)$, where under the null H_0 the means $\mu_i = 0$ identically, and under the alternative H_1 we take $\mu_i = 0$ or $\mu_i = \sqrt{2r \log n}$ for some value $r < 1$, but for which most of the μ_i are zero. Note that if $r > 1$, then the trivial test comparing $\max_i Y_i$ to $\sqrt{2 \log n}$ would be asymptotically perfect: under the null,

$$\mathbb{P}_{H_0} \left(\max_{i \leq n} Y_i \geq \sqrt{2r \log n} \right) \leq n \mathbb{P}(Y_1 \geq \sqrt{2r \log n}) \leq n \exp \left(-\frac{2r \log n}{2} \right) = n^{1-r} \rightarrow 0,$$

while under the alternative, so long as $k \gg 1$ of the signals are non-null, for $Z_i \stackrel{\text{iid}}{\sim} \mathbf{N}(0, 1)$ we have

$$\mathbb{P}_{H_1} \left(\max_{i \leq n} Y_i \geq \sqrt{2r \log n} \right) \geq \mathbb{P} \left(\max_{i \leq k} Z_i \geq 0 \right) = 1 - 2^{-k} \rightarrow 1.$$

One typical formulation is to formulate this as a mixture problem, where under H_0 we have

$$Y_i \stackrel{\text{iid}}{\sim} \mathbf{N}(0, 1)$$

and under the alternative H_1 we observe

$$Y_i \stackrel{\text{iid}}{\sim} (1 - \epsilon_n)\mathbf{N}(0, 1) + \epsilon_n\mathbf{N}(\mu, 1),$$

a mixture of $\mathbf{N}(0, 1)$ and $\mathbf{N}(\mu, 1)$ distributions, where ϵ_n determines the sparsity fraction and μ the signal strength. We then ask for the rates at which $\epsilon_n \rightarrow 0$ and the associated signal strengths $\mu > 0$ that determine whether testing between H_0 and H_1 is possible. As a brief remark, it is relatively straightforward to show that if $\epsilon_n = n^{-\beta}$ for some $\beta \in (\frac{1}{2}, 1)$, so that the signal is indeed quite sparse, then for $\mu = \sqrt{2r \log n}$, testing between H_0 and H_1 is impossible if $r < \beta - \frac{1}{2}$. (See Exercise 13.5.) \diamond

In the rest of this section, we develop some techniques to answer the questions Example 13.3.10 poses.

Abstractly, we model the sparse signal detection problem as testing mixtures, where for individual observations Y_i we have a null distribution P_0 , a known alternative P_1 , and we have a sequence of observations Y_i , $i = 1, \dots, n$, where each Y_i is drawn either

$$H_0 : Y_i \stackrel{\text{iid}}{\sim} P_0 \quad \text{or} \quad H_1 : Y_i \stackrel{\text{iid}}{\sim} (1 - \epsilon)P_0 + \epsilon P_1, \quad (13.3.15)$$

so that in the alternative H_1 we observe from a mixture of P_0 and P_1 , that is, about ϵ fraction of the time we observe data from P_1 . Then the question in such a sparse signal detection problem is the rate at which we can take $\epsilon \downarrow 0$ while still reliably testing between H_0 and H_1 .

Because the testing problem (13.3.15) is a simple hypothesis test of

$$P_0^n \quad \text{versus} \quad ((1 - \epsilon)P_0 + \epsilon P_1)^n,$$

the likelihood ratio test is always optimal, though this is a bit unsatisfying as a principal. Alternatively, by the identity that $\inf_{\Psi} \{P(\Psi = 0) + Q(\Psi = 1)\} = 1 - \|P - Q\|_{\text{TV}}$, and the equivalence between Hellinger distance and total variation distance that Proposition 2.2.7 shows, we see that

$$d_{\text{hel}}^2(P_0^n, ((1 - \epsilon)P_0 + \epsilon P_1)^n) \rightarrow 0 \quad \text{if and only if} \quad \inf_{\Psi} R_n(\Psi \mid H_0, H_1) \rightarrow 0$$

while

$$d_{\text{hel}}^2(P_0^n, ((1 - \epsilon)P_0 + \epsilon P_1)^n) \rightarrow 1 \quad \text{if and only if} \quad \inf_{\Psi} R_n(\Psi \mid H_0, H_1) \rightarrow 1.$$

Equivalently, because $d_{\text{hel}}^2(P^n, Q^n) = 1 - (1 - d_{\text{hel}}^2(P, Q))^n$, we see that for the sparse mixture testing problem (13.3.15),

$$\inf_{\Psi} R_n(\Psi \mid H_0, H_1) \rightarrow \begin{cases} 1 & \text{if } d_{\text{hel}}^2(P_0, (1 - \epsilon)P_0 + \epsilon P_1) \ll \frac{1}{n} \\ 0 & \text{if } d_{\text{hel}}^2(P_0, (1 - \epsilon)P_0 + \epsilon P_1) \gg \frac{1}{n}. \end{cases} \quad (13.3.16)$$

While a fully general theory characterizing the limits (13.3.16) does not exist, examples can help to delineate when we might hope to detect sparse signals. As a simple one that captures some of the techniques, consider the following Bernoulli detection problem.

Example 13.3.11 (Bernoulli detection): Consider a null distribution $P_0 = \text{Bernoulli}(\frac{1}{2})$ and alternatives $P_1 = \text{Bernoulli}(\frac{1+\Delta}{2})$. Noting that the mixture of Bernoulli distributions remains Bernoulli, we simply consider the Hellinger distance between P_0 and P_1 . In this case,

$$2d_{\text{hel}}^2(P_0, P_1) = 2 - \sqrt{1 + \Delta} - \sqrt{1 - \Delta} = \frac{\Delta^2}{4} + O(\Delta^3).$$

So perfect testing is asymptotically possible whenever $\Delta^2 \gg \frac{1}{n}$, and testing is asymptotically impossible whenever $\Delta^2 \ll \frac{1}{n}$. In this case, the optimal test is also simple to describe: given observations X_i , we simply test whether the sum $S_n = \sum_{i=1}^n X_i$ satisfies $S_n > \frac{n}{2} + \sqrt{\frac{n}{2} \log \frac{1}{\alpha}}$. This has level at most $\mathbb{P}_0(S_n > \frac{n}{2} + \sqrt{\frac{n}{2} \log \frac{1}{\alpha}}) \leq \exp(-\log \frac{1}{\alpha}) = \alpha$ by sub-Gaussian concentration, and its error under the alternative $X_i \stackrel{\text{iid}}{\sim} P_1$ satisfies

$$\begin{aligned} \mathbb{P}_1 \left(S_n > \frac{n}{2} + \sqrt{\frac{n}{2} \log \frac{1}{\alpha}} \right) &= \mathbb{P} \left(S_n - \frac{n}{2} - \Delta n > \sqrt{\frac{n}{2} \log \frac{1}{\alpha}} - \Delta n \right) \\ &\geq 1 - \exp \left(-\frac{2}{n} \left[\Delta n - \sqrt{\frac{n}{2} \log \frac{1}{\alpha}} \right]_+^2 \right) \rightarrow 1 \end{aligned}$$

whenever $\Delta \gg 1/\sqrt{n}$. \diamond

This example fails to capture the fuller complexity of signal detection, because no individual signal can be too strong—the observations lie in $\{0, 1\}$ regardless. Example 13.3.11 does, however, show that “interesting” regime is when the signals are quite sparse, so that $\beta > \frac{1}{2}$. Indeed, let $\beta < \frac{1}{2}$ and assume that $\|P_0 - P_1\|_{\text{TV}} \geq c > 0$. Then it is relatively simple to develop a test Ψ that achieves risk $R_n(\Psi \mid H_0, H_1) \rightarrow 0$ by counting observations more likely to have come from P_0 or P_1 , and a slight elaboration of this procedure works for $\beta = \frac{1}{2}$ as well.

Corollary 13.3.12. *Let P_0 and P_1 be a distributions satisfying $\|P_0 - P_1\|_{\text{TV}} = c > 0$, and consider the hypothesis test (13.3.15) with $\epsilon = \epsilon_n = n^{-\beta}$. Then if $\beta < \frac{1}{2}$,*

$$\inf_{\Psi} R_n(\Psi \mid H_0, H_1) \rightarrow 0.$$

If $P_{0,n}$ and $P_{1,n}$ are sequences of distributions satisfying $\|P_{0,n} - P_{1,n}\|_{\text{TV}} \rightarrow 1$, then if $\beta \leq \frac{1}{2}$,

$$\inf_{\Psi} R_n(\Psi \mid H_0, H_1) \rightarrow 0.$$

In each limit we consider the hypothesis test of $H_0 : P_{0,n}^n$ versus $H_1 : ((1 - \epsilon_n)P_{0,n} + \epsilon_n P_{1,n})^n$.

See Exercise 13.7 for a sketch of the proof.

Moving beyond signal detection problems from binary sequences, multiple hypothesis testing problems provide motivation closer to problems that arise in practice.

Example 13.3.13 (A multiple testing problem): In a multiple testing problem of n outcomes, we observe p -values, which we represent as variables $U_i \in [0, 1]$, $i = 1, \dots, n$, where $U_i \approx 0$ indicates a significant result. It is natural assume that under the null distribution, $U_i \sim \text{Uniform}([0, 1])$; more generally, we have $P(U_i \leq u) \leq u$ for $u \in [0, 1]$, meaning that under the null U_i is not particularly likely to be small. Under an alternative (i.e., something to be

discovered), U_i is likely to be nearer 0, so we have $\mathbb{P}(U_i \leq u) > u$ for $u \in [0, 1]$. We thus formulate (abstractly) the problem of detecting whether more than a negligible number of non-nulls are present as testing

$$H_0 : U_i \stackrel{\text{iid}}{\sim} \text{Uniform}([0, 1]) \quad \text{versus} \quad H_1 : U_i \stackrel{\text{iid}}{\sim} (1 - \epsilon)\text{Uniform}([0, 1]) + \epsilon P_1,$$

an instance of the abstract problem (13.3.15). This is distinct from discovering which p -values are significant and ought to be rejected—we ask whether what we observe would be unlikely if *each* p -value were null. \diamond

JCD Comment: In this example, include a version where we consider nulls $P_0 = \text{Uniform}[0, 1]$ and alternatives $P_1 = (1 - \tau)\text{Uniform}[0, \tau] + \tau\text{Uniform}[\tau, 1]$ or something similar, arguing that the testing/detection problem is the same. Could also think about it in the context of signal recovery.

Let us (heuristically) formulate Example 13.3.10 in the context of Example 13.3.13 to perform a more quantitative analysis and build a stylized problem that enables us to evaluate potential procedures for detection in Example 13.3.13. Let Φ denote the standard normal CDF, so that for detecting observations Y_i that are large, the natural p -value is $U_i = \Phi(-Y_i)$, as under the null that $Y_i \sim \mathcal{N}(0, 1)$, we have $\Phi(-Y_i) \sim \text{Uniform}([0, 1])$. Then using the approximation that $\Phi^{-1}(1 - u) \approx \sqrt{2 \log \frac{1}{u}}$ for u near 0, under the alternative $Y \sim \mathcal{N}(\mu, 1)$, for $Z \sim \mathcal{N}(0, 1)$ we have

$$\mathbb{P}(\Phi(-Y) \leq u) = \mathbb{P}(Z + \mu \geq \Phi^{-1}(1 - u) - \mu) \approx \mathbb{P}\left(Z \geq \sqrt{2 \log u^{-1}} - \mu\right).$$

If $\mu = \sqrt{2r \log n}$ as in the scaling of Example 13.3.10, then we see a transition in the probability above as $u \gtrless n^{-r}$: if $u \ll n^{-r}$, then

$$\mathbb{P}\left(Z \geq \sqrt{2 \log u^{-1}} - \sqrt{2r \log n}\right) \leq \exp\left(-(\sqrt{2 \log u^{-1}} - \sqrt{2r \log n})^2\right) \approx 0,$$

while if $u \gg n^{-r}$, then

$$\mathbb{P}\left(Z \geq \sqrt{2 \log u^{-1}} - \sqrt{2r \log n}\right) \geq 1 - \exp\left(-(\sqrt{2r \log n} - \sqrt{2 \log u^{-1}})^2\right) \approx 1.$$

That is, we have $\mathbb{P}(\Phi(-Y) \geq u) \approx 0$ if $u \gg n^{-r}$, so when Y_i comes from the non-null $\mathcal{N}(\mu, 1)$, we can (heuristically) model $\Phi(-Y_i)$ as uniform on $[0, n^{-r}]$. With this as motivation, we instantiate Example 13.3.13 with a particular choice of the alternative P_1 that exhibits this type of transitional behavior, where $U_i \sim P_1$ is uniform on $[0, \tau]$ for some $\tau > 0$, where for smaller values of τ we are more likely to observe “significant” p -values.

Example 13.3.14 (Example 13.3.13, instantiated): Consider nulls and alternatives

$$P_0 = \text{Uniform}([0, 1]) \quad \text{and} \quad P_1 = \text{Uniform}([0, \tau]), \quad \tau = n^{-r},$$

where r indicates the “signal strength”, with larger r decreasing the threshold τ . Take $\epsilon = \epsilon_n = n^{-\beta}$ for some $\beta \in (\frac{1}{2}, 1)$. Then P_0 has density 1 on $[0, 1]$, while P_1 has density $\frac{1}{\tau} \mathbf{1}_{\{0 \leq u \leq \tau\}}$. Let us compute the Hellinger distance between P_0 and the mixture, as the limit (13.3.16) dictates. The general recipe begins as follows: we let $L = \frac{dP_1}{dP_0}$ be the likelihood ratio between P_1 and P_0 , so that

$$d_{\text{hel}}^2(P_0, (1 - \epsilon)P_0 + \epsilon P_1) = \frac{1}{2} \mathbb{E}_0 \left[(1 - \sqrt{(1 - \epsilon) + \epsilon L})^2 \right] = 1 - \mathbb{E}_0 \left[\sqrt{1 - \epsilon + \epsilon L} \right],$$

and controlling these expectations gives our results.

In this example, $L = \frac{1}{\tau} \mathbf{1}\{0 \leq U \leq \tau\}$, so

$$\begin{aligned} d_{\text{hel}}^2(P_0, (1 - \epsilon)P_0 + \epsilon P_1) &= 1 - \tau \sqrt{1 - \epsilon + \epsilon/\tau} - (1 - \tau) \sqrt{1 - \epsilon} \\ &= 1 - n^{-r} \sqrt{1 - n^{-\beta} + n^{r-\beta}} - (1 - n^{-r}) \sqrt{1 - n^{-\beta}} \\ &= 1 - \sqrt{1 - n^{-\beta}} + n^{-r} \left(\sqrt{1 - n^{-\beta}} - \sqrt{1 - n^{-\beta} + n^{r-\beta}} \right) \\ &= \frac{n^{-\beta}}{2} + n^{-r} \left(\sqrt{1 - n^{-\beta}} - \sqrt{1 - n^{-\beta} + n^{r-\beta}} \right) + o(n^{-1}). \end{aligned}$$

Let d_{hel}^2 be shorthand for the squared Hellinger distance above. If $r \geq \beta$, then we have $d_{\text{hel}}^2 \gtrsim n^{-\beta} - n^{-\frac{r+\beta}{2}} \gg 1/n$, so that detection becomes trivial—some of the p -values are too small to ignore. So let us evaluate the limits when $r < \beta$. In this case,

$$\begin{aligned} \sqrt{1 - n^{-\beta}} - \sqrt{1 + n^{r-\beta} - n^{-\beta}} &= \frac{-n^{-\beta}}{2} - \frac{n^{r-\beta} - n^{-\beta}}{2} + \frac{(n^{r-\beta} - n^{-\beta})^2}{8} + O(n^{3r-3\beta}) + o(n^{-1}) \\ &= -\frac{n^{-\beta}}{2} + \frac{n^{2r-2\beta} - 2n^{r-2\beta}}{8} + O(n^{3r-3\beta}) + o(n^{-1}). \end{aligned}$$

Multiplying by n^{-r} and substituting above, we have

$$d_{\text{hel}}^2(P_0, (1 - \epsilon)P_0 + \epsilon P_1) = \frac{n^{-\beta}}{2} - \frac{n^{-\beta}}{2} + \frac{n^{r-2\beta}}{8} + O(n^{2r-3\beta}) + o(n^{-1}).$$

Finally, recognize that $2r - 3\beta < r - 2\beta$ whenever $r < \beta$, yielding the final asymptotic expression

$$d_{\text{hel}}^2(P_0, (1 - \epsilon)P_0 + \epsilon P_1) = \frac{(1 + o(1))}{8} n^{-2\beta+r} + o(n^{-1}).$$

In particular, if $r > 2\beta - 1$, then $n^{-2\beta+r} \gg 1/n$, and if $r < 2\beta - 1$, then $-3\beta + 2r < \beta - 2$. Thus, with our thresholds $\tau = n^{-r}$ and $\epsilon = n^{-\beta}$, we have

$$d_{\text{hel}}^2(P_0^n, ((1 - \epsilon)P_0 + \epsilon P_1)^n) \rightarrow \begin{cases} 1 & \text{if } r > 2\beta - 1 \\ 0 & \text{if } r < 2\beta - 1. \end{cases}$$

We see the thresholding behavior (13.3.5), where testing is possible only if the signal for the non-null p -values is suitably strong, that is, they have support $[0, n^{-r}]$ for some $r > 2\beta - 1$. \diamond

Example 13.3.14 shows that in situations in which a (relatively) small proportion of p -values in a multiple hypothesis test are from non-null distributions, so long as they concentrate near enough to 0, we can detect them. We can relax the choice to make P_0 and P_1 not absolutely continuous—they have different supports—while still providing the same asymptotic interpretation, where non-null p -values are likely to be near 0, with the same phase transitions; Exercise 13.8 provides one such modification. The development of tests that adaptively achieve the rate Example 13.3.14 suggests requires nontrivial effort. We outline one such approach in Exercise 13.6, and we also revisit Example 13.3.10 and develop its critical thresholds in the exercises as well.

JCD Comment: Finish exercises on higher criticism using the notes I've developed. Reference them here. (Write one that uses the maximum to develop, heuristically, the right threshold, then the full one.)

13.4 Instance-optimal lower bounds and super-efficiency

Super-efficiency converses allow us to confidently state that a putative benchmark for estimation performance—a lower bound—is sharp, so that outperforming the benchmark on more than a handful of problems is difficult or even impossible. Recalling Chapter 12.4, we wish to develop benchmarks that are instance specific (item (i)) and uniformly achievable (item (ii) there). To show that a benchmark or lower bound is indeed the “right” lower bound, in Chapter 12.4 we consider super-efficiency converses that hold on average, so that no procedure could outperform the lower bound except on a negligible collection of instances (recall the point (iii)). In classical estimation problems, Theorem 12.4.1 shows that the Fisher information of a parameter provides a strong indication on the limits of estimator performance. In problems of estimating or testing one-dimensional quantities and parameters of a distribution, however, we can frequently provide an alternative converse:

- (iii') There should be a pointwise *super-efficiency* converse: if a procedure outperforms the benchmark on one population P , there must exist other populations where the procedure's performance is worse than the benchmark.

Proposition 13.1.1 shows that the local modulus of continuity provides a potential benchmark lower bound, as it provides an instance-specific guarantee (desideratum (i)). The question of whether a procedure exists achieving (desideratum (ii)) the risk bound that the local modulus provides frequently depends on the problem, and unfortunately no general characterization exists (but see the bibliographic discussion for more). Nonetheless, by considering the modulus of continuity of a parameter θ with respect to the χ^2 -divergence rather than the Hellinger distance, we will obtain a benchmark for estimation complexity that, essentially, is guaranteed to always satisfy the desiderata (i) and (iii'), the pointwise super-efficiency converse.

13.4.1 Risk transfer inequalities

By using appropriate local moduli of continuity of the parameter to be estimated, we now show how to provide pointwise guarantees of the form (iii'): outperforming a benchmark at a single instance implies performing quantitatively worse at others. A technique we shall call *risk transfer inequalities* underpin such results. In some literature, these inequalities go by the name “constrained risk inequalities,” but in effort to avoid confusion with our lower bounds on constrained estimators in Chapter 11, and because the technique truly transfers “missing” risk from one problem to another, we use this nomenclature. We focus on real-valued parameters for simplicity.

Setting notation, let $\theta(P) \in \mathbb{R}$ be a parameter of interest, and for a loss $\Phi : \mathbb{R}_+ \rightarrow \mathbb{R}$, let R_P denote the risk, or expected loss,

$$R_P(\hat{\theta}) := \mathbb{E}_P \left[\Phi \left(|\hat{\theta}(X) - \theta(P)| \right) \right],$$

where $X \sim P$. We shall use the χ^2 -affinity

$$\rho_{\chi^2}(P_1 \| P_0) := D_{\chi^2}(P_1 \| P_0) + 1 = \int \frac{dP_1^2}{dP_0} = \mathbb{E}_0 \left[\frac{dP_1^2}{dP_0^2} \right] = \mathbb{E}_1 \left[\frac{dP_1}{dP_0} \right], \quad (13.4.1)$$

where \mathbb{E}_0 and \mathbb{E}_1 denote expectation under P_0 and P_1 , respectively, to show that if $R_{P_0}(\hat{\theta})$ is small under P_0 , then it must be large under alternative distributions P_1 .

The following theorem then provides a lower bound on the risk of an estimator $\hat{\theta}$ under P_1 given an upper bound on its risk under P_0 :

Theorem 13.4.1. Assume the loss $\Phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is convex. Let $\theta_0 = \theta(P_0)$ and $\theta_1 = \theta(P_1)$, and define the separation $\Delta := 2\Phi(\frac{1}{2}|\theta_0 - \theta_1|)$. If the estimator $\hat{\theta}$ satisfies $R_{P_0}(\hat{\theta}) \leq \gamma$, then

$$R_{P_1}(\hat{\theta}) \geq \left[\sqrt{\Delta} - \sqrt{\gamma \rho_{\chi^2}(P_1 \| P_0)} \right]_+^2.$$

The theorem immediately extends to product distributions, as $\rho_{\chi^2}(P_1^n \| P_0^n) = \rho_{\chi^2}(P_1 \| P_0)^n$, so that if $\hat{\theta}_n$ satisfies $\mathbb{E}_0[\Phi(|\hat{\theta}_n(X_1^n) - \theta_0|)] \leq \gamma$, then

$$\mathbb{E}_{P_1} \left[\Phi \left(|\hat{\theta}_n(X_1^n) - \theta_1| \right) \right] \geq \left[\sqrt{2\Phi(|\theta_0 - \theta_1|/2)} - \sqrt{\gamma \rho_{\chi^2}(P_1 \| P_0)^n} \right]_+^2.$$

Proof We assume without loss of generality that $\hat{\theta}(x) \in [\theta_0, \theta_1]$, as otherwise, we simply project it onto the interval $[\theta_0, \theta_1]$. For any $\theta \in [\theta_0, \theta_1]$, we have $\theta = t\theta_0 + (1-t)\theta_1$ for some $t \in [0, 1]$, so

$$\begin{aligned} \sqrt{\Phi(|\theta - \theta_0|)} + \sqrt{\Phi(|\theta - \theta_1|)} &= \sqrt{\Phi((1-t)|\theta_1 - \theta_0|)} + \sqrt{\Phi(t|\theta_1 - \theta_0|)} \\ &\geq \sqrt{\Phi((1-t)|\theta_1 - \theta_0|) + \Phi(t|\theta_1 - \theta_0|)} \geq \sqrt{2\Phi\left(\frac{1}{2}|\theta_1 - \theta_0|\right)}, \end{aligned} \quad (13.4.2)$$

because $t = \frac{1}{2}$ minimizes $\Phi(ta) + \Phi((1-t)a)$. Using the majorization (13.4.2), we thus obtain

$$\mathbb{E}_1 \left[\Phi \left(|\hat{\theta} - \theta_0| \right)^{1/2} \right] + \mathbb{E}_1 \left[\Phi \left(|\hat{\theta} - \theta_1| \right)^{1/2} \right] \geq \sqrt{2\Phi(|\theta_1 - \theta_0|/2)} = \sqrt{\Delta}.$$

Rearranging, we have by Cauchy-Schwarz that

$$R_{P_1}(\hat{\theta}) \geq \mathbb{E}_1 \left[\Phi \left(|\hat{\theta} - \theta_1| \right)^{1/2} \right]^2 \geq \left[\sqrt{\Delta} - \mathbb{E}_1 \left[\Phi \left(|\hat{\theta} - \theta_0| \right)^{1/2} \right] \right]_+^2.$$

Applying Cauchy-Schwarz once more to see that

$$\mathbb{E}_1 \left[\sqrt{\Phi(|\hat{\theta} - \theta_0|)} \right] = \mathbb{E}_0 \left[\frac{dP_1}{dP_0} \sqrt{\Phi(|\hat{\theta} - \theta_0|)} \right] \leq \mathbb{E}_0 \left[\frac{dP_1^2}{dP_0^2} \right]^{1/2} \sqrt{\mathbb{E}_0 \left[\Phi(|\hat{\theta} - \theta_0|) \right]}$$

gives the theorem. \square

Applying Theorem 13.4.1 generally involves a two step process: we assume that some estimator achieves small risk, then show that there exist distributions close in χ^2 -distance but whose parameters are different.

Example 13.4.2 (Super-efficiency in normal mean estimation): Let $\mathcal{P} = \{\mathcal{N}(\theta, 1)\}_{\theta \in \mathbb{R}}$, and assume we observe $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\theta, 1)$. Let $\hat{\theta}_n$ be an estimator satisfying

$$\mathbb{E}_{\theta_0}[|\hat{\theta}_n - \theta_0|] \leq \frac{\gamma}{\sqrt{n}}$$

at some $\theta_0 \in \mathbb{R}$. Now, consider an alternative distribution $P_1 = \mathcal{N}(\theta_1, 1)$, for which we observe that

$$\rho_{\chi^2}(P_1 \| P_0) = \exp(|\theta_1 - \theta_0|^2) \quad \text{so} \quad \rho_{\chi^2}(P_1^n \| P_0^n) = \exp(n|\theta_1 - \theta_0|^2)$$

Then by Theorem 13.4.1, we have

$$\mathbb{E}_{P_1} \left[\left| \hat{\theta}_n - \theta_1 \right| \right] \geq \left[\sqrt{|\theta_0 - \theta_1|} - \frac{1}{n^{1/4}} \exp \left(\frac{n}{2} |\theta_1 - \theta_0|^2 + \frac{n}{2} \log \gamma \right) \right]_+^2.$$

Taking $\theta_1 = \theta_0 + \sqrt{\frac{1}{n} \log \frac{1}{\gamma}}$ then yields

$$\mathbb{E}_{P_1} \left[\left| \hat{\theta}_n - \theta_1 \right| \right] \geq \frac{1}{\sqrt{n}} \left[\sqrt[4]{\log \frac{1}{\gamma}} - 1 \right]_+^2 \gtrsim \sqrt{\frac{\log \frac{1}{\gamma}}{n}}.$$

So if $\gamma \ll 1$, meaning the estimator exhibits convergence faster than $1/\sqrt{n}$ —super-efficient estimation—at θ_0 , it must pay a substantial penalty at points roughly $\sqrt{\frac{1}{n} \log \frac{1}{\gamma}}$ away. \diamond

Example 13.4.2 is an example of a broader class of super-efficiency results that arise from considering an alternative to the Hellinger modulus of continuity, where we define the χ^2 modulus, with a slight modification which does not change the asymptotic but makes computation simpler. With that in mind, recall the χ^2 -affinity (13.4.1) and define the local χ^2 -modulus of continuity by

$$\omega_{\chi^2}(\epsilon; \theta, P_0, \mathcal{P}) := \sup_{P_1 \in \mathcal{P}} \{ |\theta(P_1) - \theta(P_0)| \mid \log \rho_{\chi^2}(P_1 \| P_0) \leq 2\epsilon^2 \}. \quad (13.4.3)$$

By Proposition 2.2.9, $2d_{\text{hel}}^2(P_0, P_1) \leq \log(1 + D_{\chi^2}(P_1 \| P_0)) = \log \rho_{\chi^2}(P_1 \| P_0)$, so the Hellinger modulus dominates the χ^2 -modulus:

$$\omega_{\chi^2}(\epsilon; \theta, P_0, \mathcal{P}) \leq \omega_{\text{hel}}(\epsilon; \theta, P_0, \mathcal{P}).$$

As we have seen, however, in finite-dimensional problems, where divergences between distributions are locally quadratic (13.1.7) as in Examples 13.1.7 and 13.1.8, the moduli are equivalent to within constant factors. Proposition 13.1.4, coupled with a guarantee that the χ^2 -divergence behaves locally quadratically (as in Lemma 13.1.10) shows this.

Corollary 13.4.3. *Let the conditions of Theorem 13.4.1 hold, and define the shorthand $\omega(\epsilon) = \omega_{\chi^2}(\epsilon; \theta, P_0, \mathcal{P})$. Let $\hat{\theta}_n$ be any estimator satisfying*

$$\mathbb{E}_{P_0} \left[\Phi(|\hat{\theta}_n(X_1^n) - \theta_0|) \right] \leq \gamma \cdot \Phi(\omega(1/\sqrt{n}))$$

for some $\gamma < 1$. Then for all distributions P_1 such that $\log(1 + D_{\chi^2}(P_1 \| P_0)) \leq \frac{1}{n} \log \frac{1}{\gamma}$,

$$\mathbb{E}_{P_1} \left[\Phi \left(|\hat{\theta}_n(X_1^n) - \theta| \right) \right] \geq \left[\sqrt{2\Phi(|\theta(P_1) - \theta_0|/2)} - \sqrt{\Phi(\omega(1/\sqrt{n}))} \right]_+^2.$$

In particular,

$$\sup_{\log \rho_{\chi^2}(P_1 \| P_0) \leq \frac{1}{n} \log \frac{1}{\gamma}} \mathbb{E}_{P_1} \left[\Phi \left(|\hat{\theta}_n(X_1^n) - \theta| \right) \right] \geq \left[\sqrt{2\Phi \left(\frac{1}{2} \omega \left(\left(\frac{1}{n} \log \frac{1}{\gamma} \right)^{1/2} \right) \right)} - \sqrt{\Phi(\omega(1/\sqrt{n}))} \right]_+^2.$$

Proof Let $\Phi_n = \Phi(\omega_{\chi^2}(\frac{1}{\sqrt{n}}; \theta, P_0, \mathcal{P}))$ for shorthand. By Theorem 13.4.1, for any P_1 and associated parameters θ_0, θ_1 , we have

$$R_{P_1}(\hat{\theta}) \geq \left[\sqrt{2\Phi(|\theta_1 - \theta_0|/2)} - \sqrt{\gamma\Phi_n \cdot \rho_{\chi^2}(P_1\|P_0)^n} \right]_+^2.$$

Take P_1 to be any distribution satisfying $\log \rho_{\chi^2}(P_1\|P_0) \leq \frac{1}{n} \log \frac{1}{\gamma}$, so that $\gamma\Phi_n \cdot \rho_{\chi^2}(P_1\|P_0)^n \leq \gamma\Phi_n \exp(\log \frac{1}{\gamma}) = \Phi_n$. Take a supremum over such P_1 . \square

So we see that, if an estimator outperforms the risk the χ^2 -modulus of continuity predicts at a distribution P_0 , then its loss on nearby distributions P_1 must be larger than that the modulus predicts. A few subtleties are worth noting here, however: the moduli in Corollary 13.4.3 are local to P_0 , not to the alternative P_1 , meaning that we have not fully satisfied the pointwise super-efficiency converse (iii'). Of course, if the modulus (13.4.3) itself is appropriately continuous in the argument P_0 , then we obtain a super-efficiency result.

Example 13.4.4 (Superefficiency in regular parametric families): Let $\{P_\theta\}_{\theta \in \mathbb{R}^d}$ be a parametric family regular enough for the χ^2 -divergence (Definition 13.1), so that $D_{\chi^2}(P_{\theta+v}\|P_\theta) = v^\top J(\theta)v + o(\|v\|^2)$ for v small, where $J(\theta) := \mathbb{E}_\theta[\dot{\ell}_\theta \dot{\ell}_\theta^\top]$ is the Fisher information matrix for the parameter θ . Let $T : \mathbb{R}^d \rightarrow \mathbb{R}$ be a continuously differentiable function of interest and \hat{T}_n be an estimator of $T(\theta)$ satisfying

$$\mathbb{E}_{\theta_0} \left[(\hat{T}_n(X_1^n) - T(\theta_0))^2 \right] \leq \gamma \frac{\nabla T(\theta_0)^\top J(\theta_0)^{-1} \nabla T(\theta_0)}{n}$$

at some θ_0 . We recognize $\nabla T(\theta)^\top J(\theta)^{-1} \nabla T(\theta)$ as the constant factors characterizing the local modulus of continuity for problems with Fisher information, as in Proposition 13.1.4. Then applying Corollary 13.4.3 with the squared error $\Phi(t) = t^2$, we obtain that for all θ satisfying $(\theta - \theta_0)^\top J(\theta_0)(\theta - \theta_0) \lesssim \frac{1}{n} \log \frac{1}{\gamma}$ that

$$\mathbb{E}_\theta \left[(\hat{T}_n(X_1^n) - T(\theta))^2 \right] \geq \left[\sqrt{\frac{|T(\theta) - T(\theta_0)|^2}{n}} - O(1) \frac{\|J(\theta_0)^{-1/2} \nabla T(\theta_0)\|_2}{\sqrt{n}} \right]_+^2.$$

Using the differentiability of T , we have $|T(\theta) - T(\theta_0)| = |\langle \nabla T(\theta_0), \theta - \theta_0 \rangle| + o(\|\theta - \theta_0\|)$. Then choosing

$$\theta = c \sqrt{\frac{1}{n} \log \frac{1}{\gamma}} \frac{J(\theta_0)^{-1/2} \nabla T(\theta_0)}{\|J(\theta_0)^{-1/2} \nabla T(\theta_0)\|_2}$$

for some numerical constant $c > 0$ to see that

$$\mathbb{E}_\theta \left[(\hat{T}_n(X_1^n) - T(\theta))^2 \right] \gtrsim \frac{\nabla T(\theta_0)^\top J(\theta_0)^{-1} \nabla T(\theta_0)}{n} \log \frac{1}{\gamma} \gtrsim \frac{\nabla T(\theta)^\top J(\theta)^{-1} \nabla T(\theta)}{n} \log \frac{1}{\gamma},$$

where we use the continuity of J and $\nabla T(\theta)$. \diamond

Example 13.4.4 extends Example 13.4.2, showing that estimating too accurately at any one point implies much worse estimation performance elsewhere.

13.4.2 A general risk transfer bound

The convexity of the loss in Theorem 13.4.1 may be unsatisfying, as in some cases we wish to lower bound quantities such as the probability of error, which corresponds to the indicator of an estimate belonging to a set or not, and as such certainly lacks convexity. Therefore, we extend Theorem 13.4.1 to apply to general loss functions.

Corollary 13.4.5. *Let the loss $\Phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ be non-decreasing. Let $\theta_0 = \theta(P_0)$ and $\theta_1 = \theta(P_1)$, and define the separation $\Delta := \Phi(\frac{1}{2}|\theta_0 - \theta_1|)$. If the estimator $\hat{\theta}$ satisfies $R_{P_0}(\hat{\theta}) \leq \gamma$, then*

$$R_{P_1}(\hat{\theta}) \geq \left[\sqrt{\Delta} - \sqrt{\gamma \rho_{\chi^2}(P_1 \| P_0)} \right]_+^2.$$

Proof The only modification we need to make to the proof of Theorem 13.4.1 is to replace the majorization inequality (13.4.2) with a slightly smaller lower bound. To that end, let $\theta \in [\theta_0, \theta_1]$ so that $\theta = t\theta_0 + (1-t)\theta_1$ for some $t \in [0, 1]$. Then at least one of $|\theta - \theta_0| \geq \frac{1}{2}|\theta_1 - \theta_0|$ or $|\theta - \theta_1| \geq \frac{1}{2}|\theta_1 - \theta_0|$, and so

$$\sqrt{\Phi(|\theta - \theta_0|)} + \sqrt{\Phi(|\theta - \theta_1|)} \geq \sqrt{\Phi(|\theta_1 - \theta_0|/2)}.$$

The rest of the proof is identical. \square

Applying Corollary 13.4.5 follows with much the same technique as our applications of Theorem 13.4.1 have. To highlight the techniques, we revisit normal mean estimation, Example 13.4.2, but now consider a testing scenario.

Example 13.4.6 (Super-efficient testing of a normal mean): Let $\mathcal{P} = \{N(\theta, \sigma^2)\}_{\theta \in \mathbb{R}}$ be the collection of normal distributions, and consider a zero-one loss function indicating whether the estimated mean is near the true mean, $\Phi(t) = \mathbf{1}\{|t| \geq \sigma/\sqrt{n}\}$, so that

$$R_{P_\theta^n}(\hat{\theta}) = P_\theta^n \left(|\hat{\theta}(X_1^n) - \theta| \geq \frac{\sigma}{\sqrt{n}} \right).$$

We expect that the risk of most estimators should be roughly of constant order, as minimax considerations dictate. Let us assume however, that we have a sequence of estimators $\hat{\theta}_n$ satisfying $P_0^n(|\hat{\theta}_n| \geq \frac{\sigma}{\sqrt{n}}) \leq \delta_n$, where $\delta_n \rightarrow 0$ (i.e., $\hat{\theta}_n$ is super-efficient at $\theta = 0$). Fix any $0 < c < 1$, and define the sequence of local alternative parameter spaces

$$\Theta_n := \left\{ \theta \in \mathbb{R} \mid 2\frac{\sigma}{\sqrt{n}} \leq |\theta| \leq \frac{\sigma}{\sqrt{n}} \sqrt{c \log \frac{1}{\delta_n}} \right\}.$$

Then we claim that

$$\liminf_n \inf_{\theta \in \Theta_n} P_\theta^n \left(|\hat{\theta}_n(X_1^n) - \theta| \geq \frac{\sigma}{\sqrt{n}} \right) = 1, \quad (13.4.4)$$

that is, for a reasonably large collection of parameters θ in a shell around 0, the estimator is *never* within σ/\sqrt{n} of the true parameter.

To see the limit (13.4.4), assume that n is large enough that $c \log \frac{1}{\delta_n} \geq 2$, so that Θ_n is non-empty. Let $\theta \in \Theta_n$. Then calculating the χ^2 -affinity, we obtain

$$\rho_{\chi^2}(P_\theta \| P_0)^n = \exp \left(\frac{n\theta^2}{\sigma^2} \right) \leq \exp \left(\frac{c\sigma^2 n \log \frac{1}{\delta_n}}{\sigma^2 n} \right) = \delta_n^{-c}.$$

For any $\theta \in \Theta_n$, we have $\Phi(\frac{1}{2}|\theta|) = \mathbf{1}\{|\theta| \geq \sigma/\sqrt{n}\} = 1$, so Corollary 13.4.5 gives

$$R_{P_\theta^n}(\hat{\theta}_n) \geq \left[1 - \delta_n^{(1-c)/2}\right]_+^2 \rightarrow 1,$$

implying inequality (13.4.4). \diamond

13.4.3 Risk transfer with mixtures

Section 13.2 develops lower bound techniques via Le Cam's convex hull method, which can provide stronger minimax lower bounds than the two-point methods, which work with hardest 1-dimensional subproblems. Sometimes, as in Section 13.2.3, the two-point lower bounds are loose even in the rate of convergence. Nonetheless, an extended (local) modulus of continuity with respect to the χ^2 -divergence over convex hulls of distributions still admits risk transfer inequalities, making it possible to extend the pointwise super-efficiency converses of Section 13.4.1 to more sophisticated problems. For simplicity in this section, and because the extension is more or less trivial given the development thus far, we focus only on absolute losses $\Phi(t) = |t|$.

To that end, we define the *mixture modulus*, which extends the local modulus of continuity with respect to the χ^2 -divergence (13.4.3), by

$$\bar{\omega}_{\chi^2}(\epsilon; \theta, P_0, \mathcal{P}) := \sup_{m \in \mathbb{N}} \sup_{P_1, \dots, P_m \in \mathcal{P}} \left\{ \min_{i \leq m} \{|\theta(P_0) - \theta(P_i)|\} \mid \log \rho_{\chi^2}(\bar{P} \| P_0) \leq 2\epsilon^2 \text{ for some } \bar{P} \in \text{Conv}\{P_i\}_{i=1}^m \right\}. \quad (13.4.5)$$

The quantity (13.4.5) is more sophisticated than the typical local modulus, but measures how much it is possible to perturb the parameter θ of interest while making sure that some element of the convex hull $\text{Conv}\{P_i\}_{i=1}^m$ is close to the base distribution P_0 . We abuse notation slightly to let $\mathcal{P}^n = \{P^n\}_{P \in \mathcal{P}}$ be the collection of product distributions and

$$\bar{\omega}_{\chi^2}(\epsilon; \theta, P_0^n, \mathcal{P}^n) := \sup_{m \in \mathbb{N}} \sup_{P_1, \dots, P_m \in \mathcal{P}} \left\{ \min_{i \leq m} \{|\theta(P_0) - \theta(P_i)|\} \mid \log \rho_{\chi^2}(\bar{P}^n \| P_0^n) \leq 2\epsilon^2 \text{ for some } \bar{P}^n \in \text{Conv}\{P_i^n\}_{i=1}^m \right\}.$$

The mixture modulus will provide a benchmark sufficient to guarantee pointwise super-efficiency converses. First, however, we note that as a corollary of the convex hull method in Theorem 13.2.1, that the mixture modulus always lower bounds the minimax risk.

Corollary 13.4.7. *Let the conditions of Theorem 13.2.1 hold. Then for each n and $P_0 \in \mathcal{P}$,*

$$\mathfrak{M}_n(\theta(\mathcal{P}), |\cdot|) \geq \frac{1}{4} \bar{\omega}_{\chi^2}\left(\frac{1}{2}; \theta, P_0^n, \mathcal{P}^n\right).$$

Proof Let $P_0 \in \mathcal{P}$ and $\mathcal{P}_1 \subset \mathcal{P}$ be any δ -separated collection, that is, satisfying $|\theta(P_0) - \theta(P)| \geq \delta$ for all $P \in \mathcal{P}_1$. Then Theorem 13.2.1 implies that

$$\mathfrak{M}_n(\theta(\mathcal{P}), |\cdot|) \geq \frac{\delta}{2} (1 - \|P_0^n - \bar{P}^n\|_{\text{TV}})$$

for any $\bar{P}^n \in \text{Conv}\{\mathcal{P}_1^n\}$. By Pinsker's inequality (Proposition 2.2.8) and Proposition 2.2.9,

$$2 \|P_0^n - \bar{P}^n\|_{\text{TV}}^2 \leq \log \rho_{\chi^2}(\bar{P}^n \| P_0^n),$$

so if $\log \rho_{\chi^2}(\overline{P^n} \| P_0^n) \leq \frac{1}{2}$ then $\|P_0^n - \overline{P^n}\|_{TV} \leq \frac{1}{2}$. Taking a supremum over the convex hull gives the result. \square

We turn to risk transfer inequalities, where we show in analogy with Theorem 13.4.1 that achieving small risk at some distribution P_0 implies achieving larger risk at others.

Theorem 13.4.8. *Let $P_v \in \mathcal{P}$ for $v = 0, 1, \dots, M$, and let $\theta_v = \theta(P_v)$ for each v , and define the separation $\Delta = \min_{v \geq 1} |\theta_0 - \theta_v|$. If the estimator $\hat{\theta}$ satisfies $\mathbb{E}_{P_0}[|\hat{\theta} - \theta_0|] \leq \gamma$, then for any nonnegative $\lambda \in \mathbb{R}_+^M$ with $\langle \lambda, \mathbf{1} \rangle = 1$ and $P_\lambda := \sum_{v=1}^M \lambda_v P_v$,*

$$\sum_{v=1}^M \lambda_v \mathbb{E}_{P_v} [|\hat{\theta} - \theta_v|] \geq \left[\sqrt{\Delta} - \sqrt{\gamma \rho_{\chi^2}(P_\lambda \| P_0)} \right]_+^2.$$

Proof The proof mirrors that of Theorem 13.4.1. Fix $v \in [M]$. Then for any $\theta \in \mathbb{R}$, we have as in inequality (13.4.2) that

$$\sqrt{|\theta - \theta_0|} + \sqrt{|\theta - \theta_v|} \geq \sqrt{|\theta_0 - \theta_v|} \geq \sqrt{\Delta}.$$

So we see that for any v ,

$$\mathbb{E}_v [|\hat{\theta} - \theta_v|^{1/2} + |\hat{\theta} - \theta_0|^{1/2}] \geq \sqrt{\Delta},$$

implying by Cauchy-Schwarz that

$$\mathbb{E}_v [|\hat{\theta} - \theta_v|] \geq \mathbb{E}_v [|\hat{\theta} - \theta_v|^{1/2}]^2 \geq \left[\sqrt{\Delta} - \mathbb{E}_v [|\hat{\theta} - \theta_0|^{1/2}] \right]_+^2.$$

As λ forms a convex combination, the convexity of $t \mapsto [t]_+^2$ yields by Jensen's inequality that

$$\sum_v \lambda_v \mathbb{E}_v [|\hat{\theta} - \theta_v|] \geq \left[\sqrt{\Delta} - \sum_v \lambda_v \mathbb{E}_v [|\hat{\theta} - \theta_0|^{1/2}] \right]_+^2.$$

Applying the Cauchy-Schwarz inequality once again, we obtain

$$\begin{aligned} \sum_v \lambda_v \mathbb{E}_v [|\hat{\theta} - \theta_0|^{1/2}] &= \mathbb{E}_0 \left[\frac{dP_\lambda}{dP_0} |\hat{\theta} - \theta_0|^{1/2} \right] \leq \mathbb{E}_0 \left[\frac{dP_\lambda^2}{dP_0^2} \right]^{1/2} \mathbb{E}_0 [|\hat{\theta} - \theta_0|]^{1/2} \\ &= \rho_{\chi^2}(P_\lambda \| P_0)^{1/2} \mathbb{E}_0 [|\hat{\theta} - \theta_0|]^{1/2}. \end{aligned}$$

The theorem follows. \square

We can leverage Theorem 13.4.8 to show how the convex hull modulus (13.4.5) implies super-efficiency converses. For shorthand in the corollary, we define

$$\bar{\omega}_n(\epsilon) := \bar{\omega}_{\chi^2}(\epsilon; \theta, P_0^n, \mathcal{P}^n).$$

Corollary 13.4.9. *Let $P_0 \in \mathcal{P}$, and let $\hat{\theta}_n$ be an estimator satisfying $\mathbb{E}_{P_0}[|\hat{\theta}_n(X_1^n) - \theta(P_0)|] \leq \gamma \bar{\omega}_n(1)$. Then there exists a collection of distributions $\{P_v\}_{v=1}^M \subset \mathcal{P}$ and a vector $\lambda \in \mathbb{R}_+^M$, $\langle \mathbf{1}, \lambda \rangle = 1$, such that*

$$\sum_{v=1}^M \lambda_v \mathbb{E}_{P_v^n} [|\hat{\theta}_n(X_1^n) - \theta_v|] \geq \left[\sqrt{\bar{\omega}_n(1/\sqrt{\gamma})} - \sqrt{\bar{\omega}_n(1)} \right]_+^2.$$

Proof By Theorem 13.4.8, for any collection $P_1, \dots, P_M \in \mathcal{P}$ with associated parameters $\theta_v = \theta(P_v)$, if we let $\Delta = \min_v |\theta_v - \theta_0|$ then

$$\sum_{v=1}^M \lambda_v \mathbb{E}_{P_v^n} \left[|\hat{\theta}(X_1^n) - \theta_v| \right] \geq \left[\sqrt{\Delta} - \sqrt{\gamma \bar{\omega}_n(1) \rho_{\chi^2}(P_\lambda^n \| P_0^n)} \right]_+^2$$

for any vector $\lambda \in \mathbb{R}_+^M$ with $\langle \mathbf{1}, \lambda \rangle = 1$, where $P_\lambda^n = \sum_v \lambda_v P_v^n$. So long as $\log \rho_{\chi^2}(P_\lambda^n \| P_0^n) \leq \frac{1}{\gamma}$, we therefore have

$$\sum_{v=1}^M \lambda_v \mathbb{E}_{P_v^n} \left[|\hat{\theta}(X_1^n) - \theta_v| \right] \geq \left[\sqrt{\Delta} - \sqrt{\bar{\omega}_n(1)} \right]_+^2.$$

Take a supremum over all such convex combinations of $\mathcal{P}^n = \{P^n\}_{P \in \mathcal{P}}$ and let $\Delta = \bar{\omega}_n(1/\sqrt{\gamma})$. \square

Corollary 13.4.9 shows how the convex-hull modulus around P_0 thus provides (roughly) an unimprovable bound: if an estimator achieves a faster convergence rate by a factor $\gamma < 1$ than $\bar{\omega}_n(1)$, then on average over some collection of alternative distributions, the risk of the estimator must scale at least as the modulus across a larger neighborhood.

Example 13.4.10 (An integral functional bound): We revisit the estimation problems in Section 13.2.3, where the goal was to estimate $T_k(f) := \int_0^1 (f^{(k)}(x))^2 dx$ given observations $Y_i = f(X_i) + \varepsilon_i$. For simplicity, let us take $k = s = 1$, so that f is assumed to have bounded second derivative and we wish to estimate $T_1(f) = \int f'(x)^2 dx$. Suppose that we have an estimator that, for the identically 0 function $f_0(x) = 0$, has risk

$$\mathbb{E}_{f_0} [|\hat{T}_n|] \leq \gamma_n n^{-\frac{4}{9}}, \quad \text{where } \gamma_n \rightarrow 0,$$

which is (asymptotically) faster than the lower bound in Proposition 13.2.8. Then Corollary 13.4.9 implies that for the class $\mathcal{F}_2 = \{f : [0, 1] \rightarrow \mathbb{R} \mid \|f''\|_\infty \leq 1, f(0) = 0\}$, we have

$$\frac{\sup_{f \in \mathcal{F}_2} \mathbb{E}_f [|\hat{T}(X_1^n, Y_1^n) - T_1(f)|]}{n^{-4/9}} \rightarrow \infty.$$

To see this limit, we revisit the construction in the proof of Proposition 13.2.8. Letting P_f be the induced distribution of observations under $Y_i = f(X_i) + \varepsilon_i$, where P_0 is $Y_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$, equations (13.2.7) and (13.2.8) imply the following: there exist positive numerical constants c_0, c_1, c_2 such that for each $m, n \in \mathbb{N}$, there is a collection of functions \mathcal{G} such that

$$T_1(f) \geq \frac{c_0}{m^2} \text{ for } f \in \mathcal{G} \quad \text{and} \quad \rho_{\chi^2}(\bar{P}^n \| P_0) \leq \exp \left(\frac{c_1 n^2}{m^9} + \frac{nc_2}{m^8} \right),$$

where $\bar{P}^n = \frac{1}{|\mathcal{G}|} \sum_{f \in \mathcal{G}} P_f^n$. In particular, this implies that

$$\bar{\omega}_{\chi^2}(\epsilon; T_1, P_0^n, \mathcal{P}^n) \gtrsim \sup_{m \in \mathbb{N}} \left\{ \frac{1}{m^2} \mid \frac{n^2}{m^9} + \frac{n}{m^8} \leq \epsilon^2 \right\} \gtrsim \left(\frac{\epsilon}{n} \right)^{4/9},$$

where we chose $m = (n/\epsilon)^{2/9}$. By Corollary 13.4.9, we thus have

$$\sup_{f \in \mathcal{F}_2} \mathbb{E}_{P_f^n} [|\hat{T}_n - T_1(f)|] \gtrsim \left(\frac{1}{n\gamma_n} \right)^{4/9} \gg \frac{1}{n^{4/9}},$$

a quantitative lower bound on the risk of \hat{T}_n . \diamond

13.5 Deferred and technical proofs

13.5.1 Proof of Lemma 13.1.6

We mostly reduce to the 1-dimensional case by considering changing $v_s \rightarrow v \in \mathbb{S}^{d-1}$ as $s \rightarrow 0$, because the claim in the lemma is equivalent to the claim that

$$\limsup_{v \rightarrow 0} \frac{1}{\|v\|^2} \left(d_{\text{hel}}^2(P_{\theta_0+v}, P_{\theta_0}) - \frac{1}{8} v^\top J(\theta_0) v \right) = 0.$$

For $s \in \mathbb{R}$ define $f_s(x) := p_{\theta_0+sv_s}(x)$, so that $\dot{f}_s(x) = \langle \dot{p}_{\theta_0+sv_s}(x), v \rangle$ and the score $\dot{\ell}_s(x) := \frac{\dot{f}_s(x)}{f_s(x)} = \langle v_s, \dot{p}_{\theta_0+sv_s}(x) \rangle / p_{\theta_0+sv_s}(x)$. Now for $t > 0$,

$$\sqrt{f_t(x)} = \sqrt{f_0(x)} + \int_0^t \frac{\dot{f}_s(x)}{2\sqrt{f_s(x)}} ds,$$

so we have

$$\frac{1}{t^2} d_{\text{hel}}^2(P_{\theta_0+vt}, P_{\theta_0}) = \frac{1}{8} \int \left(\frac{1}{t} \int_0^t \frac{\dot{f}_s(x)}{\sqrt{f_s(x)}} ds \right)^2 d\mu(x).$$

By Jensen's inequality, we have the domination condition

$$\left(\frac{1}{t} \int_0^t \frac{\dot{f}_s(x)}{\sqrt{f_s(x)}} ds \right)^2 \leq \frac{1}{t} \int_0^t \frac{\dot{f}_s(x)^2}{f_s(x)} ds = \frac{1}{t} \int_0^t \dot{\ell}_s^2(x) f_s(x) ds,$$

and by Fubini's theorem this quantity is integrable:

$$\int \left(\frac{1}{t} \int_0^t \frac{\dot{f}_s(x)}{\sqrt{f_s(x)}} ds \right)^2 d\mu(x) \leq \frac{1}{t} \int \int_0^t \dot{\ell}_s^2(x) f_s(x) ds d\mu(x) = \frac{1}{t} \int_0^t v_s^\top J(\theta_0 + sv_s) v_s ds < \infty$$

for small enough t , as J is assumed continuous at θ_0 . The assumption that for μ -almost all x there is a neighborhood of θ_0 on which $\sqrt{p_\theta}$ is continuously differentiable then allows us to apply a variant of dominated convergence (see Proposition A.3.3 in Appendix A.3), because

$$\left(\frac{1}{t} \int_0^t \frac{\langle \dot{p}_{\theta_0+sv_s}(x), v_s \rangle}{\sqrt{p_{\theta_0+sv_s}(x)}} ds \right)^2 \rightarrow \frac{\langle v, \dot{p}_{\theta_0}(x) \rangle^2}{p_{\theta_0}(x)}$$

for μ -almost all x . Then $\frac{1}{t^2} d_{\text{hel}}^2(P_{\theta_0+tv_t}, P_{\theta_0}) \rightarrow \frac{1}{8} \int \langle \dot{\ell}_{\theta_0}(x), v \rangle^2 p_{\theta_0}(x) d\mu(x) = \frac{1}{8} v^\top J(\theta_0) v$. Because we allow $v_s \rightarrow v$ to vary, this implies the desired uniformity.

13.5.2 Proof of Lemma 13.1.10

We prove the expansions related to the χ^2 -divergence first, showing that

$$\int \left(\frac{p_{\theta+v}(x)}{p_\theta(x)} - 1 - \dot{\ell}_\theta(x)^\top v \right)^2 p_\theta(x) d\mu(x) = o(\|v\|^2). \quad (13.5.1)$$

Assuming the bound (13.5.1), standard L^p -convergence give the expansion $D_{\chi^2}(P_{\theta+v} \| P_\theta) = v^\top J(\theta)v + o(\|v\|^2)$, as $\mathbb{E}_\theta[\dot{\ell}_\theta \dot{\ell}_\theta^\top] = J(\theta)$. As in the proof of Lemma 13.1.6, we essentially reduce to the 1-dimensional case. Let $v_t \rightarrow v \in \mathbb{S}^{d-1}$ as $t \rightarrow 0$. Then

$$\frac{1}{t^2} \int \left(\frac{dP_{\theta+tv_t}}{dP_\theta} - 1 - t\dot{\ell}_\theta^\top v_t \right)^2 dP_\theta \leq \int h^2(x) \|v_t\|^2 dP_\theta(x) \lesssim \int h^2(x) dP_\theta(x) < \infty,$$

and using that for μ -almost all x

$$\frac{1}{t} \left(\frac{p_{\theta+tv_t}(x)}{p_\theta(x)} - 1 - \dot{\ell}_\theta(x)^\top v_t \right) \rightarrow 0$$

as $t \rightarrow 0$, we have by dominated convergence that

$$\int \left(\frac{dP_{\theta+tv_t}}{dP_\theta} - 1 - t\dot{\ell}_\theta^\top v_t \right)^2 dP_\theta = o(t^2),$$

giving the claim (13.5.1).

Now we demonstrate the expansions of the f -divergences. Without loss of generality, we assume that $f'(1) = 0$, because we may always replace f with $\tilde{f}(t) := f(t) - tf'(1) + f'(1)$, which satisfies $\tilde{f}(1) = \tilde{f}'(1) = 0$ and gives the same divergence. Once again let $v_t \rightarrow v \in \mathbb{S}^{d-1}$ as $t \rightarrow 0$, and assume w.l.o.g. that $\|v_t\|_2 \leq 2$ for all t . Define the densities $q_t(x) = p_{\theta+tv_t}(x)$, so that by assumption

$$\left| \frac{q_t(x)}{q_0(x)} - 1 - t\dot{\ell}_\theta(x)^\top v_t \right| \leq 2th(x).$$

Define the gap

$$g_t(x) := \frac{1}{t^2} \left| f\left(\frac{q_t(x)}{q_0(x)}\right) - \frac{f''(1)}{2} \left(\frac{q_t(x)}{q_0(x)} - 1\right)^2 \right|.$$

Then

$$\begin{aligned} g_t(x) &\leq \frac{C}{t^2} \left(\frac{q_t(x)}{q_0(x)} - 1 \right)^2 \leq \frac{2C}{t^2} \left(\frac{q_t(x)}{q_0(x)} - 1 - t\dot{\ell}_\theta(x)^\top v_t \right)^2 + \frac{2C}{t^2} t^2 \left\| \dot{\ell}_\theta(x) \right\|_2^2 \|v_t\|_2^2 \\ &\leq 2C \left(h(x)^2 + \left\| \dot{\ell}_\theta(x) \right\|_2^2 \right) \|v_t\|_2^2 \leq 4C \left(h(x)^2 + \left\| \dot{\ell}_\theta(x) \right\|_2^2 \right), \end{aligned}$$

which is integrable with respect to P_θ . Because $f(1+s) = f(1) + f'(1)s + \frac{f''(1)}{2}s^2 + o(s^2)$ and $f(1) = f'(1) = 0$ by our w.l.o.g. assumption,

$$\begin{aligned} g_t(x) &= \frac{1}{t^2} \left| f\left(1 + \frac{q_t(x)}{q_0(x)} - 1\right) - \frac{f''(1)}{2} \left(\frac{q_t(x)}{q_0(x)} - 1\right)^2 \right| \\ &= o\left(\frac{1}{t^2} \left(\frac{q_t(x)}{q_0(x)} - 1\right)^2\right) = o\left(t^{-2} \left(t\dot{\ell}_\theta(x)^\top v_t + o(t)\right)^2\right) \rightarrow 0 \end{aligned}$$

for all x . Dominated convergence then implies that $\lim_{t \rightarrow 0} \int g_t(x) q_0(x) d\mu(x) = 0$, and so

$$\frac{1}{t^2} \left(D_f(Q_t \| Q_0) - \frac{f''(1)}{2} D_{\chi^2}(Q_t \| Q_0) \right) \rightarrow 0$$

as $t \rightarrow 0$. Because we took $q_t(x) = p_{\theta+tv_t}$ along a path, this gives the result.

13.6 Bibliography

JCD Comment: We stole the mixture idea from David Pollard I believe. Might want to cite that old Tsybakov paper for constrained risk inequalities.

Cite Donoho and Liu, Geometrizing Rates II, for the attainability and add some commentary on it. Mention that there are extensions to privacy as well.

Cite Cai and Low for the motivation of super-efficiency converses, but also recognize that they've been around for a long time and cite Hodges counterexample, and cite van der Vaart's lecture on super-efficiency.

Cite Birgé and Massart [33] for smooth functionals.

Outline

- I. Motivation: function values, testing certain quantities (e.g. is $\|P - Q\|_{\text{TV}} \geq \epsilon$ or not), entropy and other quantities, and allows superefficiency guarantees in an elegant way
- II. Le Cam's methods
 1. The general form with mixtures
 2. The χ^2 -type bounds, with mixtures to a point mass
 3. Geometrizing rates of convergence
 4. Examples: Fisher information in classical problems (especially for a one-dimensional quantity)
 5. Example: testing distance to uniformity (failure from standard two-point bound)
 6. More sophisticated examples:
 - a. Smooth functionals (as in Birgé and Massart [33]), like differential entropy $\int h(x) \log h(x) dx$
 - b. Higher-dimensional problems, which are hard
- III. "Best possible" lower bounds, super-efficiency and constrained risk inequalities
 1. Basic (two-point) constrained risk inequality (cf. [71])
 2. Constrained risk inequality when P_1 is actually a mixture (easiest with a functional): means that any minimax bound around P_0 is quite strong
 3. Potentially (?): Cai and Low [44] paper on minimax estimation for $\frac{1}{n} \|\theta\|_1$ when $y = \theta + \varepsilon$ in a Gaussian sequence model as an example and application of a constrained risk inequality. This is probably too challenging, though—can we find a case where polynomials actually allow us to do stuff?
 - a. Hard because of all the polynomial approximation stuff... but maybe there is a simpler version that simply shows how approximation via polynomials allows lower bounds. Approach works for Gaussian stuff, as in Cai and Low [44] or the earlier paper "Effect of mean on variance function estimation in nonparametric regression" by Wang, Brown, Cai, Levine.
 - b. Similar idea gives variation distance bounds for Poisson priors on parameters when seeking lower bounds on estimating entropy $H(X) = -\sum_x p_x \log p_x$ of discrete distributions with (unknown) support; see [191].

13.7 A useful divergence calculation

JCD Comment: Put this as exercises, related to Examples in Section 13.2.3.

Now, let us suppose that we define the collection $\{P_v\}$ by tiltings of an underlying base distribution P_0 , where each tilting is indexed by a function $g_v : \mathcal{X} \rightarrow [-1, \infty)$, and where

$$dP_v(x) = (1 + g_v(x))dP_0(x),$$

while $\int g_v dP_0 = 0$, so that each P_v is a valid distribution. Let P_v^n be the distribution of n observations $X_i \stackrel{\text{iid}}{\sim} P_v$, and let $\overline{P^n} = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P_v^n$.

Lemma 13.7.1. *Define the inner product $\langle f, g \rangle_P = \int f(x)g(x)dP(x)$ and let $V, V' \stackrel{\text{iid}}{\sim} \text{Uniform}(\mathcal{V})$. Then*

$$D_{\chi^2}(\overline{P^n} \| P_0) + 1 \leq \mathbb{E}[\exp(n\langle g_V, g_{V'} \rangle_{P_0})].$$

Proof The simple technical lemma 13.2.3 essentially gives us the result. We observe that

$$D_{\chi^2}(\overline{P^n} \| P_0^n) + 1 = \frac{1}{|\mathcal{V}|^2} \sum_{v, v'} \int \frac{dP_v^n dP_{v'}^n}{dP_0^n} = \frac{1}{|\mathcal{V}|^2} \sum_{v, v'} \left(\int (1 + g_v(x))(1 + g_{v'}(x))dP_0(x) \right)^n$$

because $P_v^n(x_1, \dots, x_n) = \prod_{i=1}^n (1 + g_v(x_i))dP_0(x_i)$, so that the integral decomposes into a product of integrals. Then expanding $(1 + g_v)(1 + g_{v'})$ and noting that each has zero mean under P_0 gives

$$D_{\chi^2}(\overline{P^n} \| P_0^n) + 1 = \frac{1}{|\mathcal{V}|^2} \sum_{v, v'} (1 + \mathbb{E}_0[g_v(X)g_{v'}(X)])^n.$$

Lastly, we note that $(1 + t) \leq e^t$ for all t , and so

$$\frac{1}{|\mathcal{V}|^2} \sum_{v, v'} (1 + \mathbb{E}_0[g_v(X)g_{v'}(X)])^n \leq \frac{1}{|\mathcal{V}|^2} \sum_{v, v'} \exp(n\mathbb{E}_0[g_v(X)g_{v'}(X)]),$$

which is of course equivalent to the result we desired. \square

A specialization of Lemma 13.7.1 follows when we choose our functions g to correspond to a partition of \mathcal{X} -space. Here, we define the following.

Definition 13.2. *Let $k \in \mathbb{N}$ and the functions $\phi_j : \mathcal{X} \rightarrow [-b, b]$. Then the functions ϕ_j are an admissible partition with variances σ_j^2 of \mathcal{X} with respect to a probability distribution P_0 if*

- (i) *The supports $E_j = \text{supp } \phi_j$ of each of the functions are disjoint.*
- (ii) *Each function has P_0 mean 0, i.e., $\mathbb{E}_{P_0}[\phi_j(X)] = 0$ for each j .*
- (iii) *Function j has variance $\sigma_j^2 = \mathbb{E}_{P_0}[\phi_j^2(X)] = \int \phi_j^2(x)dP_0(x)$.*

With such a partition, we can define the functions $g_v(x) = t\langle v, \phi(x) \rangle = t \sum_{j=1}^k v_j \phi_j(x)$ for $|t| \leq 1/b$, and if we take $\mathcal{V} = \{-1, 1\}^k$, we obtain the following.

Lemma 13.7.2. *Let the functions $\{\phi_j\}_{j=1}^k$ be an admissible partition of \mathcal{X} with variances σ_j^2 . Fix $|t| \leq \frac{1}{b}$, and let $dP_{tv} = (1 + t\langle v, \phi(x) \rangle) dP_0(x)$ and $\bar{P}_t^n = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P_v^n$. Then*

$$D_{\chi^2}(\bar{P}_t^n \| P_0) \leq \exp\left(\frac{n^2 t^4}{2} \sum_{j=1}^k \sigma_j^4\right) - 1,$$

and if $|t| \leq \frac{1}{\sqrt{n} (\sum_{j=1}^k \sigma_j^4)^{1/4}}$, then

$$D_{\chi^2}(\bar{P}_t^n \| P_0) \leq n^2 t^4 \sum_{j=1}^k \sigma_j^4.$$

Proof First, if $\phi(x) = [\phi_j(x)]_{j=1}^k$, then $\mathbb{E}_0[\phi(X)\phi(X)^T] = \text{diag}(\sigma_j^2)$, that is, the diagonal matrix with σ_j^2 on its diagonal. By Lemma 13.7.1, we therefore have

$$D_{\chi^2}(\bar{P}_t^n \| P_0) + 1 \leq \mathbb{E} \left[\exp \left(nt^2 \sum_{j=1}^k \sigma_j^2 V_j V_j' \right) \right] \leq \mathbb{E} \left[\exp \left(\frac{n^2 t^4}{2} \sum_{j=1}^k \sigma_j^4 V_j^2 \right) \right]$$

by Hoeffding's Lemma (see Example 4.1.6), as $V_j \stackrel{\text{iid}}{\sim} \text{Uniform}(\{\pm 1\})$. Noting that $V_j^2 = 1$ gives the first part of the lemma. The final statement is immediate once we observe that $e^x \leq 1 + (e-1)x \leq 1 + 2x$ for $0 \leq x \leq 1$. \square

13.8 Exercises

Exercise 13.1: Recall the Hellinger distance between distributions P and Q with densities p, q is $d_{\text{hel}}(P, Q)^2 = \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx$. Let P be $\mathcal{N}(\mu_0, \Sigma)$ and Q be $\mathcal{N}(\mu_1, \Sigma)$. Show that

$$\frac{1}{2} d_{\text{hel}}(P, Q)^2 = 1 - \exp \left(-\frac{1}{8} (\mu_0 - \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1) \right).$$

Exercise 13.2: Demonstrate the claims in Example 13.1.9.

Exercise 13.3 (The Hodges super-efficient estimator): Consider a normal mean estimation problem, $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\theta, 1)$ where $\theta \in \mathbb{R}$ is unknown, and define the *Hodges estimator*

$$\hat{\theta}_H(X_1^n) := \begin{cases} 0 & \text{if } |\bar{X}_n| \leq n^{-1/4} \\ \bar{X}_n & \text{otherwise.} \end{cases}$$

(a) Show that the limiting distribution of $\hat{\theta}_H$ satisfies

$$\sqrt{n}(\hat{\theta}_H(X_1^n) - \theta) \stackrel{d}{\rightsquigarrow} \begin{cases} \mathbf{1}_0 & \text{if } \theta = 0 \\ \mathcal{N}(0, 1) & \text{otherwise.} \end{cases}$$

That is, in a pointwise sense, the asymptotic distribution of the Hodges estimator is the same of the sample mean, except at 0 where it is simply the point mass at 0.

(b) Show that for any $c < 1$,

$$\lim_{n \rightarrow \infty} \inf_{\theta} \left\{ P_{\theta} \left(|\hat{\theta}_H(X_1^n) - \theta| \geq 1/\sqrt{n} \right) \mid \frac{2}{\sqrt{n}} \leq |\theta| \leq \frac{c}{\sqrt{2}} \frac{1}{n^{1/4}} \right\} = 1.$$

What does this say about the performance of $\hat{\theta}_H$?

Exercise 13.4: Suppose that the test Ψ has test risk for testing between \mathcal{P}_0 and \mathcal{P}_1 satisfying $R_n(\Psi \mid \mathcal{P}_0, \mathcal{P}_1) \leq \frac{1}{3}$. Let $k \in \mathbb{N}$. Show how, given a sample of size kn , we can develop a test Ψ^* with

$$R_{kn}(\Psi^* \mid \mathcal{P}_0, \mathcal{P}_1) \leq 2 \exp(-ck),$$

where $c > 0$ is a numerical constant. *Hint.* Split the sample into k samples of size n , and then apply Ψ to each.

Exercise 13.5: Take the sampling model (13.3.15), where $\epsilon = \epsilon_n = n^{-\beta}$ for some $\beta \in (\frac{1}{2}, 1)$, and $P_1 = \mathcal{N}(\mu_n, 1)$ for $\mu_n = \sqrt{2r \log n}$.

(a) Show that

$$1 + D_{\chi^2}((1 - \epsilon)P_0 + \epsilon P_1 \parallel P_0) = 1 - \epsilon^2 + \epsilon^2 e^{\mu^2}.$$

(b) Show that if $r < \beta - \frac{1}{2}$, then

$$d_{\text{hel}}^2(P_0, (1 - \epsilon)P_0 + \epsilon P_1) \ll 1/n.$$

Hint. Use Proposition 2.2.9.

(c) Conclude that if $\mu_n \leq \sqrt{2r \log n}$ for $r < \beta - \frac{1}{2}$, it is asymptotically impossible to test between $H_0 : P_0^n$ and $H_1 : ((1 - \epsilon)P_0 + \epsilon P_1)^n$.

Exercise 13.6 (Higher criticism and sparse detection [68]): One version of Tukey's *Higher Criticism* statistic is to consider the normalized process

$$Z_n(t) := \sqrt{n} \frac{P_n(U \leq t) - t}{\sqrt{t(1-t)}} \quad \text{and} \quad T_n^{\text{HC}} := \sup_t \{Z_n(t) \mid 1/n \leq t \leq 1 - 1/n\}.$$

This higher criticism statistic converges in probability: $T_n^{\text{HC}} / \sqrt{2 \log \log n} \xrightarrow{P} 1$ when $U_i \stackrel{\text{iid}}{\sim} \text{Uniform}([0, 1])$ (see, e.g., [169, Theorem 16.1.2]). Consider the null and alternative setting of Example 13.3.14.

(a) Show that for any $\alpha \in (0, 1)$, there is a sequence $a_n \rightarrow 1$ such that under the null H_0 ,

$$\limsup_n \mathbb{P}_{H_0} \left(T_n^{\text{HC}} \geq a_n \sqrt{2 \log \log n} \right) \leq \alpha.$$

(b) Show that under the alternative H_1 that $U_i \stackrel{\text{iid}}{\sim} (1 - \epsilon)\text{Uniform}[0, 1] + \epsilon\text{Uniform}[0, \tau]$, where $\epsilon = n^{-\beta}$ and $\tau = n^{-r}$ for some $\beta \in (\frac{1}{2}, 1)$ and $r > 2\beta - 1$, we have

$$\mathbb{P}_{H_1} \left(T_n^{\text{HC}} \geq a_n \sqrt{2 \log \log n} \right) \rightarrow 1.$$

Thus, higher criticism has optimal power to detect anywhere on the interior of the region $r > 2\beta - 1$.

Exercise 13.7 (Non-sparse detection problems are easy): Let P_0 and P_1 be distributions on a random variable X and consider observations X_i drawn i.i.d. from

$$H_0 : X_i \stackrel{\text{iid}}{\sim} P_0 \quad \text{or} \quad H_1 : X_i \stackrel{\text{iid}}{\sim} (1 - \epsilon)P_0 + \epsilon P_1.$$

We investigate the rate at which we may allow $\epsilon = \epsilon_n \rightarrow 0$ while still testing accurately, taking $\epsilon_n = n^{-\beta}$ for some $\beta \in [0, 1]$.

- (a) Let $\|P_0 - P_1\|_{\text{TV}} \geq c > 0$, and assume that $\beta < \frac{1}{2}$. Let A be such that $P_0(A) + P_1(A^c) \leq 1 - c$. Construct a test Ψ_n based on $S_n := \sum_{i=1}^n \mathbf{1}\{X_i \in A\}$ that achieves $R_n(\Psi_n \mid H_0, H_1) \rightarrow 0$ as $n \rightarrow \infty$.
- (b) Assume that as $n \rightarrow \infty$, we allow P_1 to vary so that $\|P_0 - P_1\|_{\text{TV}} \rightarrow 1$, and $\beta \leq \frac{1}{2}$. Show that it is possible to construct a test Ψ_n such that $R_n(\Psi_n \mid H_0, H_1) \rightarrow 0$ as $n \rightarrow \infty$.
- (c) Let $P_0 = \text{N}(0, 1)$ and $P_1 = \text{N}(\mu, 1)$ for $\mu = \sqrt{2r \log n}$, where $r > 0$. Show that $\|P_0 - P_1\|_{\text{TV}} \geq 1 - n^{-r}$. What does this say about Example 13.3.10?
- (d) Let $P_0 = \text{Uniform}([0, 1])$ and $P_1 = \text{Uniform}([0, \tau])$ for any $\tau < 1$. Give $\|P_0 - P_1\|_{\text{TV}}$. What does this say about Example 13.3.13?

Exercise 13.8 (A less brutal multiple testing scenario): Instead of the instantiation of Example 13.3.13 in Example 13.3.14, consider nulls and alternatives

$$P_0 = \text{Uniform}([0, 1]) \quad \text{and} \quad P_1 = (1 - \tau)\text{Uniform}([0, \tau]) + \tau\text{Uniform}([\tau, 1]),$$

where again $\epsilon = n^{-\beta}$ while $\epsilon \ll \tau \ll 1$. (So $U \sim P_1$ has density $\frac{1-\tau}{\tau} \mathbf{1}\{0 \leq u \leq \tau\} + \frac{\tau}{1-\tau} \mathbf{1}\{\tau \leq u \leq 1\}$, making P_1 and P_0 absolutely continuous.)

- (a) Show that for $\epsilon \ll \tau$,

$$d_{\text{hel}}^2(P_0, (1 - \epsilon)P_0 + \epsilon P_1) = \frac{1 + o(1)}{8} \frac{\epsilon^2}{\tau}.$$

- (b) Do the asymptotics in the conclusion of Example 13.3.14 still hold in this case? Why or why not?

Exercise 13.9 (Multiple testing and the Benjamini-Hochberg procedure): Consider the Benjamini-Hochberg step-up procedure (9.8.3) for rejecting hypotheses in Exercise 9.5. Here you provide a few results that suggest a failure mode of this procedure. (Though of course it exhibits many complementary optimality properties.)

- (a) Let $U_i \stackrel{\text{iid}}{\sim} \text{Uniform}[0, 1]$. For $c \in [0, 1/n]$, define the region

$$A_c := \{x \in \mathbb{R}^n \mid 0 \leq x_1 \leq \dots \leq x_n \leq 1, x_i \geq ci \text{ for each } i\}.$$

Letting $U_{(1)} \leq U_{(2)} \leq \dots \leq U_{(n)}$ be the order statistics of $U \in [0, 1]^n$, show that

$$\mathbb{P}((U_{(1)}, \dots, U_{(n)}) \in A_c) = \mathbb{P}(U_{(1)} \geq c, U_{(2)} \geq 2c, \dots, U_{(n)} \geq nc) = 1 - nc.$$

Hint. The density of the collection of order statistics is uniform on A_0 and satisfies $p(u_1^n) = 1/n!$.

- (b) Let \hat{k} be the number of hypothesis the BH procedure (9.8.3) rejects on data $U_i \stackrel{\text{iid}}{\sim} \text{Uniform}[0, 1]$. Show that $\mathbb{P}(\hat{k} > 0) = \alpha$.
- (c) Suppose that $U_1, \dots, U_k \stackrel{\text{iid}}{\sim} [0, \tau]$, where $k = n^{1-\beta}$ and $\tau = n^{-r}$, mimicking the setting of Example 13.3.14, where $2\beta - 1 < r < \beta$. Show that if we only compute order statistics for this set of variables,

$$\mathbb{P}\left(\text{any } U_{(i)} \geq \frac{\alpha i}{n}\right) \rightarrow 0.$$

Exercise 13.10 (A phase diagram): Consider the multiple hypothesis testing (sparse p -values) scenario of Example 13.3.14 or its less brutal version in Exercise 13.8. Use the results of the example and Exercise 9.5 to justify the phase diagram in Figure 13.1, which graphically depicts the following:

- i. If $r > \beta$, then it is possible to asymptotically perfectly recover all null and non-null signals; if $r < \beta$, it is impossible.
- ii. If $r > 2\beta - 1$, then it is possible to detect that there are non-null signals, but not to identify them; if $r < 2\beta - 1$, then even detection is impossible.

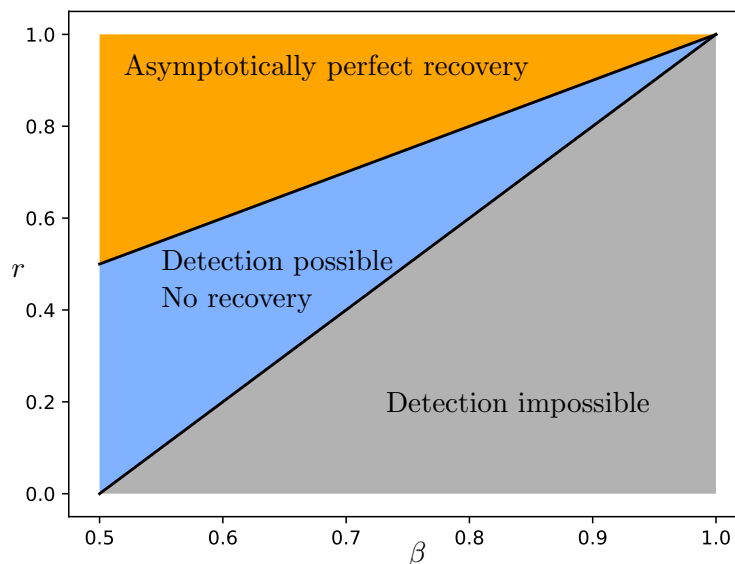


Figure 13.1. Phase diagram for the stylized sparse hypothesis testing in Example 13.3.14. See Exercise 13.10.

JCD Comment: Exercise ideas:

1. Work through lower bound for sparse testing. Also add a reference to the next problem in the main text.

Exercise 13.11 (Detecting a sparse signal through the maximum):

- (a) Show that if $Z_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$, then $\max_{i \leq n} Z_i / \sqrt{2 \log n} \xrightarrow{p} 1$. *Hint.* Use Mills ratio, Exercise 4.3.

- (b) Show that if $Z_i \stackrel{\text{iid}}{\sim} \mathbf{N}(0, 1)$, then for any $\alpha \in [0, 1]$, there exists a sequence $a_n = a_n(\alpha) \rightarrow 1$ for which $\mathbb{P}(\max_{i \leq n} Z_i \geq a_n \sqrt{2 \log n}) \rightarrow \alpha$.

Recall the testing problem of Example 13.3.10 to distinguish $H_0 : Y_i \stackrel{\text{iid}}{\sim} \mathbf{N}(0, 1)$ from the alternative $H_1 : Y_i \stackrel{\text{iid}}{\sim} (1 - \epsilon_n)\mathbf{N}(0, 1) + \epsilon_n\mathbf{N}(\mu_n, 1)$, where $\epsilon_n = n^{-\beta}$ and $\mu_n = \sqrt{2r \log n}$, and $\beta \in (\frac{1}{2}, 1)$.

- (c) Show that if $r > (1 - \sqrt{1 - \beta})^2$, then for $a_n = a_n(0) \rightarrow 1$ from part (b), the test

$$\Psi_n := \begin{cases} 0 & \text{if } \max_{i \leq n} Y_i \leq a_n \sqrt{2 \log n} \\ 1 & \text{if } \max_{i \leq n} Y_i > a_n \sqrt{2 \log n} \end{cases}$$

satisfies $R_n(\Psi_n | H_0, H_1) \rightarrow 0$ if $r > (1 - \sqrt{1 - \beta})^2$.

Exercise 13.12 (Poissonization: lower bounds [191]): Prove the lower bound in Proposition 13.3.7, inequality (13.3.9), that is, that for numerical constants C, c ,

$$\mathfrak{M}_{\text{Poi}(2n)} - Cr^2 \exp(-cn) \leq \mathfrak{M}_n.$$

Hint. Bound $\mathfrak{M}_{\text{Poi}(2n)}$ with a weighted sum of \mathfrak{M}_m . Use the MGF calculation that for $X \sim \text{Poi}(\lambda)$, $\mathbb{E}[e^{tX}] = \exp(\lambda(e^t - 1))$ to show that $N \sim \text{Poi}(2n)$ is concentrated above n .

Exercise 13.13 (Poissonization: upper bounds [191]): Assume the minimax result that

$$\mathfrak{M}_n = \sup_{\pi} \inf_{T_n} \mathbb{E}[(T_n(X_1^n) - T(p))^2],$$

where the supremum is over probability distributions (priors π) on $p \in \Delta_k$, and the expectation is now over the random choice of p and the sample $X_1^n \stackrel{\text{iid}}{\sim} p$ drawn conditional on p . (This is a standard infinite-dimensional saddle point result generalizing von-Neumann's minimax theorem; cf. [87, 170].) You will show the upper bound in Proposition 13.3.7, Eq. (13.3.9).

Let $\{T_m\}$ be an arbitrary sequence of estimators and define the sequence of averaged risks

$$r_m := \mathbb{E}[(T_m(X_1^m) - T(p))^2].$$

Define the modified risks $\tilde{r}_m = \min\{r_1, \dots, r_m\} = \min\{\tilde{r}_{m-1}, r_m\}$, and the “corrected” estimators

$$\tilde{T}_m(x_1^m) := \begin{cases} \tilde{T}_{m-1}(x_1^{m-1}) & \text{if } r_m \geq \tilde{r}_{m-1} \\ T_m(x_1^m) & \text{if } r_m < \tilde{r}_{m-1}. \end{cases}$$

- (a) Show that $\mathbb{E}[(\tilde{T}_m(X_1^m) - T(p))^2] \leq \mathbb{E}[(T_m(X_1^m) - T(p))^2]$.

- (b) Show that

$$\frac{1}{2} \inf_{T_n} \mathbb{E}[(T_n(X_1^n) - T(p))^2] \leq \mathbb{E}[(T_N(X_1^N) - T(p))^2]$$

for $N \sim \text{Poi}(n/2)$ and $p \sim \pi$, then X_i drawn i.i.d. conditionally on p .

- (c) Finalize the proof of the upper bound in inequality (13.3.9).

Exercise 13.14: Consider the hypothesis testing problem of testing whether a collection of independent Bernoulli random variables X_1, \dots, X_n is fair (H_0 , so that $\mathbb{P}(X_i = 1) = \frac{1}{2}$ for each i) or that there are unfair subcollections. That is, we wish to test

$$\begin{aligned} H_0 : X_i &\stackrel{\text{iid}}{\sim} \text{Bernoulli}(\tfrac{1}{2}) \\ H_1 : X_i &\stackrel{\text{ind}}{\sim} \text{Bernoulli}(\tfrac{1+\theta_i}{2}), \theta \in C \end{aligned}$$

for a set $C \subset [-1, 1]^n$. Show that if the set C is orthosymmetric, meaning that whenever $\theta \in C$ then $S\theta \in C$ for any diagonal matrix S of signs, i.e. $\text{diag}(S) \in \{\pm 1\}^n$, then no test can reliably distinguish H_0 from H_1 (in a minimax sense). *Hint.* Let $v \in \mathcal{V} := \{\pm 1\}^n$ index coordinate signs and define $\theta_v = Dv$ for some diagonal D , where $Dv \in C$. Let P_v be the product distribution with $X_i \sim \text{Bernoulli}(\frac{1+D_i v_i}{2})$. What is $\frac{1}{2^n} \sum_{v \in \mathcal{V}} P_v$?

Exercise 13.15 (Testing a trend in independent Bernoullis): Consider testing whether a collection of Bernoulli random variables has an “upward trend” over time, by which we mean that if $X_i \sim \text{Bernoulli}(p_i)$ independently, then

$$\bar{p}_{\text{end}} := \frac{1}{n/4} \sum_{i=\frac{3n}{4}+1}^n p_i > \bar{p}_{\text{beg}} := \frac{1}{n/4} \sum_{i=1}^{n/4} p_i.$$

Consider the following more quantitative version of this problem: we wish to test

$$\begin{aligned} H_0 : X_i &\stackrel{\text{iid}}{\sim} \text{Bernoulli}(\tfrac{1}{2}) \\ H_1 : X_i &\stackrel{\text{ind}}{\sim} \text{Bernoulli}(p_i), \bar{p}_{\text{end}} - \bar{p}_{\text{beg}} \geq \delta. \end{aligned}$$

(a) Use Le Cam’s two-point method to show that there exists a numerical constant $c > 0$ such that for $\delta \leq \frac{c}{\sqrt{n}}$, no test can reliably distinguish H_0 from H_1 .

(b) Use the statistic

$$T_n := \frac{1}{n/4} \sum_{i=\frac{3n}{4}+1}^n X_i - \frac{1}{n/4} \sum_{i=1}^{n/4} X_i$$

to develop a test Ψ (use Proposition 13.3.2) that achieves test risk $R_n(\Psi \mid H_0, H_1) \leq \frac{1}{4}$ whenever $\delta \geq \frac{C}{\sqrt{n}}$, where $C < \infty$ is a constant.

Exercise 13.16: Prove the identity (13.3.11).

Exercise 13.17 (Unbiased estimators of distance for multinomials): Let $X_i \stackrel{\text{iid}}{\sim} p$, $i = 1, \dots, n$, and $Y_i \stackrel{\text{iid}}{\sim} q$, $i = 1, \dots, m$, meaning that X_1^n and Y_1^m are multinomial samples for $p, q \in \Delta_d$. Define the empirical estimators $\hat{p}_j = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i = j\}$ and $\hat{q}_j = \frac{1}{m} \sum_{i=1}^m \mathbf{1}\{Y_i = j\}$.

(a) Give $\mathbb{E}[\|\hat{p}\|_2^2]$.

(b) Show that $T_n := \|\hat{p} - \hat{q}\|_2^2$ satisfies

$$\mathbb{E}[T_n] = \|p - q\|_2^2 + \frac{1}{n} + \frac{1}{m} - \frac{1}{n} \|p\|_2^2 - \frac{1}{m} \|q\|_2^2.$$

(c) Modify T_n into a new statistic T_n^{unb} so that $\mathbb{E}[T_n^{\text{unb}}] = \|p - q\|_2^2$.

Exercise 13.18: Show that in the hypothesis testing problem (13.3.13), there is a numerical constant $c > 0$ such that $\delta \leq c/\sqrt{n}$ implies that no test can reliably distinguish H_0 from H_1 .

Exercise 13.19: Consider the linear regression model

$$Y = X\theta + \varepsilon, \quad \varepsilon \sim \mathbf{N}(0, I_n)$$

where the design $X \in \mathbb{R}^{n \times d}$ is fixed (and known), and $Y \in \mathbb{R}^n$ and $\theta \in \mathbb{R}^d$. Consider the two hypotheses

$$H_0 : \theta = 0 \quad \text{versus} \quad H_1 : \|\theta\|_2 \geq r.$$

Use the convex hull method for lower bounds in testing (Proposition 13.3.1) to show that if

$$r \leq \frac{d^{1/4}}{\sqrt{n} \|n^{-1/2} X\|_{\text{op}}}$$

then any test of H_0 against H_1 has minimax test risk at least $\frac{1}{2}$. If X is an orthogonal design, so that $\frac{1}{n} X^\top X = I_d$, is this result tight?

JCD Comment: Put that additional question in here from the Scratch directory

JCD Comment:

1. Poissonization: remark in main text.
2. Work through Liam's ℓ_1 -multinomial testing
3. Lower bound for testing whether collection of coins is fair or some number are unfair.

Part III

Entropy, predictions, divergences, and information

JCD Comment: TODO: include an intro describing how we're going to get into operational interpretations of these quantities, and that is fun.

Chapter 14

Predictions, loss functions, and entropies

In prediction problems broadly construed, we have a random variable X and a label, or target or response, Y , and we wish to fit a model or predictive function that accurately predicts the value of Y given X . There are several perspectives possible when we consider such problems, each with attendant advantages and challenges. We can roughly divide these into three approaches, though there is considerable overlap between the tools, techniques, and goals of the three:

- (1) Point prediction, where we wish to find a prediction function f so that $f(X)$ most accurately predicts Y itself.
- (2) Probabilistic prediction, where we output a predicted distribution P of Y , and we seek $\mathbb{P}(Y = y \mid X = x) \approx P(Y = y \mid X = x)$, where here \mathbb{P} denotes the “true” probability and P the predicted one. A relaxed version of this is *calibration*, the subject of the next chapter, where we ask that $\mathbb{P}(Y = y \mid P) \approx P(Y = y)$, that is, the distribution of Y given a predicted distribution P is accurate.
- (3) Predictive inference, where for a given level $\alpha \in (0, 1)$, we seek a confidence set mapping C such that $\mathbb{P}(Y \in C(X)) \approx 1 - \alpha$.

We focus mostly on the former two, though there is overlap between the approaches.

In this first chapter of the sequence, we focus on the probabilistic prediction problem. Our main goal will be to elucidate and identify loss functions for choosing probabilistic predictions that are *proper*, meaning that the true distribution of Y minimizes the loss, and *strictly proper*, meaning that the true distribution of Y uniquely minimizes the loss. As part of this, we will develop mappings between losses and entropy-type functionals; these will repose on convex analytic techniques for their cleanest statements, highlighting the links between convex analysis, prediction, and information. Moreover, we highlight how *any* proper loss (which will be defined) is in correspondence with a particular measure of entropy on the distribution P , and how these connect with an object known as the *Bregman divergence* central to convex optimization. For the deepest understanding of this chapter, it will therefore be useful to review the basic concepts of convexity (e.g., convex sets, functions, and subgradients) in Appendix B, as well as the more subtle tools on optimality and stability of solutions to convex optimization problems in Appendix C. We give an overview of the important results in Section 14.1.1.

14.1 Proper losses, scoring rules, and generalized entropies

As motivation, consider a weather forecasting problem: a meteorologist wishes to prediction the weather Y_t on days $t = 1, 2, \dots$, where $Y_t = 1$ indicates rain and $Y_t = 0$ indicates no rain. At time t , using covariates X_t (for example, the weather the previous day, long term trends, or simulations), the forecaster predicts a probability $p_t \in [0, 1]$. We would like the forecaster's predictions to be as accurate as possible, so that $\mathbb{P}(Y_t = 1) \approx p_t$. Following the standard dicta of decision theory, we choose a loss function $\ell(p, y)$ that scores a prediction p for a given outcome y . Ideally, the forecaster should have an incentive to make predictions as accurately as possible, so the distribution minimizing the expected loss should coincide with the true distribution of Y .

This leads to proper losses. In our treatment, we will sometimes allow infinite values, so we work with the upper and lower extended real lines, recalling that $\bar{\mathbb{R}} = \mathbb{R} \cup \{+\infty\}$ and $\underline{\mathbb{R}} = \mathbb{R} \cup \{-\infty\}$.

Definition 14.1. Let \mathcal{P} be a collection of distributions on \mathcal{Y} . A loss $\ell : \mathcal{P} \times \mathcal{Y} \rightarrow \bar{\mathbb{R}}$ is proper if, whenever $Y \sim P \in \mathcal{P}$,

$$\mathbb{E}[\ell(P, Y)] \leq \mathbb{E}[\ell(Q, Y)] \text{ for all } Q \in \mathcal{P}.$$

The loss is strictly proper if the preceding inequality is strict whenever $Q \neq P$.

In much of the literature on prediction, one instead considers *proper scoring rules*, which are simply negative proper losses, that is, functions $S : \mathcal{P} \times \mathcal{Y}$ satisfying $S(P, y) = -\ell(P, y)$ for a (strictly) proper loss. We focus on losses for consistency with the convex analytic tools we develop. In addition, frequently we will work with discrete distributions, so that Y has a probability mass function (p.m.f.), in which case we will use $p \in \Delta_k := \{p \in \mathbb{R}_+^k \mid \langle \mathbf{1}, p \rangle = 1\}$ to identify the distribution and $\ell(p, y)$ instead of $\ell(P, y)$.

Perhaps the two most famous proper losses are the log loss and the squared loss (often termed *Brier scoring*). For simplicity let us assume that $\mathcal{Y} \in \{1, 2, \dots, k\}$, and let $\Delta_k = \{p \in \mathbb{R}_+^k \mid \mathbf{1}^T p = 1\}$ be the probability simplex; we then identify distributions P on \mathcal{Y} with vectors $p \in \Delta_k$, and abuse notation to write $\ell(p, y)$ accordingly and when it is unambiguous. The squared loss is then

$$\ell_{\text{sq}}(p, y) = (p_y - 1)^2 + \sum_{i \neq y} p_i^2 = \|p - e_y\|_2^2,$$

where e_y is the y th standard basis vector, while the log loss (really, the negative logarithm) is

$$\ell_{\log}(p, y) = -\log p_y.$$

Both of these are strictly proper. To this propriety, let Y have p.m.f. $p \in \Delta_k$, so that $\mathbb{P}(Y = y) = p_y$. Then for the squared loss and any $q \in \Delta_k$, we have

$$\mathbb{E}[\ell_{\text{sq}}(q, Y)] - \mathbb{E}[\ell_{\text{sq}}(p, Y)] = \mathbb{E}[\|q - e_Y\|_2^2] - \mathbb{E}[\|p - e_Y\|_2^2] = \|q\|_2^2 - 2\langle q, p \rangle + 2\langle p, p \rangle = \|q - p\|_2^2.$$

For the log loss, we have

$$\mathbb{E}[\ell_{\log}(q, Y)] - \mathbb{E}[\ell_{\log}(p, Y)] = -\sum_{y=1}^k p_y \log q_y + \sum_{y=1}^k p_y \log p_y = \sum_{y=1}^k p_y \log \frac{p_y}{q_y} = D_{\text{kl}}(p \| q).$$

It is immediate that $q = p$ uniquely minimizes each loss.

That the gap between the expected losses at q and p reduced to a particular divergence-like measure—the squared ℓ_2 -distance in the case of the squared loss and the KL-divergence in the

case of the log loss—is no accident. In fact, for proper losses, we will show that this divergence representation necessarily holds.

The key underlying our development is a particular construction, which we present in Section 14.1.2, that transforms a loss into a generalized notion of entropy. Because it is so central, we highlight it here, though before doing so, we take a brief detour through a few of the concepts in convexity we require. Figures representing these results capture most of the mathematical content, while Chapters B and C in the appendices contain proofs of the results we require.

14.1.1 A convexity primer

Recall that a function $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ is convex if for all $x, y \in \text{dom } f$ and $\lambda \in [0, 1]$, we have

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y),$$

where for $x \notin \text{dom } f$ we define $f(x) = +\infty$. We exclusively work with proper convex functions, so that $f(x) > -\infty$ for each x . Typically, we work with closed convex f , meaning that the epigraph $\text{epi } f = \{(x, t) \in \mathbb{R}^d \times \mathbb{R} \mid f(x) \leq t\} \subset \mathbb{R}^{d+1}$ is a closed set; equivalently, f is lower semi-continuous, so that $\liminf_{y \rightarrow x} f(y) \geq f(x)$. A concave function f is one for which $-f$ is convex.

Three main concepts form the basis for our development. The first is the *subgradient* (see Appendix B.3). For a function $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$, the *subgradient set* (also called the subdifferential) at the point x is

$$\partial f(x) := \left\{ s \in \mathbb{R}^d \mid f(y) \geq f(x) + \langle s, y - x \rangle \text{ for all } y \in \mathbb{R}^d \right\}. \quad (14.1.1)$$

If f is a convex function, then at any point x in the relative interior of its domain, $\partial f(x)$ is non-empty (Theorem B.3.3). Moreover, a quick calculation shows that x minimizes $f(x)$ if and only if $0 \in \partial f(x)$, and (a more challenging calculation) that if $\partial f(x) = \{s\}$ is a singleton, then $s = \nabla f(x)$ is the usual gradient. See the left plot of Figure 14.1. We shall in some cases allow subgradients to take values in the extended reals $\overline{\mathbb{R}}^k$ and $\underline{\mathbb{R}}^k$, which will necessitate some additional care.

The second concept is that the supremum of a collection of convex functions is always convex, that is, if f_α is convex for each index $\alpha \in \mathcal{A}$, then

$$f(x) := \sup_{\alpha \in \mathcal{A}} f_\alpha(x)$$

is convex, and f is closed in f_α is closed for each α . The closure of f is immediate because $\text{epi } f = \bigcap \text{epi } f_\alpha$, and convexity follows because

$$f(\lambda x + (1 - \lambda)y) \leq \sup_{\alpha \in \mathcal{A}} \{\lambda f_\alpha(x) + (1 - \lambda)f_\alpha(y)\} \leq \lambda \sup_{\alpha \in \mathcal{A}} f_\alpha(x) + (1 - \lambda) \sup_{\alpha \in \mathcal{A}} f_\alpha(y).$$

Conveniently, subdifferentiability of individual f_α implies the subdifferentiability of f when the supremum is attained. Indeed, let $\mathcal{A}(x) = \{\alpha \mid f_\alpha(x) = f(x)\}$. Then

$$\partial f(x) \subset \text{Conv} \{s_\alpha \in \partial f_\alpha(x) \mid \alpha \in \mathcal{A}(x)\} \quad (14.1.2)$$

because if $s = \sum_{\alpha \in \mathcal{A}(x)} \lambda_\alpha s_\alpha$ for some $\lambda_\alpha \geq 0$ with $\sum_\alpha \lambda_\alpha = 1$, then

$$f(y) \geq \sum_{\alpha \in \mathcal{A}(x)} \lambda_\alpha f_\alpha(y) \geq \sum_{\alpha \in \mathcal{A}(x)} \lambda_\alpha [f_\alpha(x) + \langle s_\alpha, y - x \rangle] = f(x) + \langle s, y - x \rangle.$$

See the right plot of Figure 14.1.

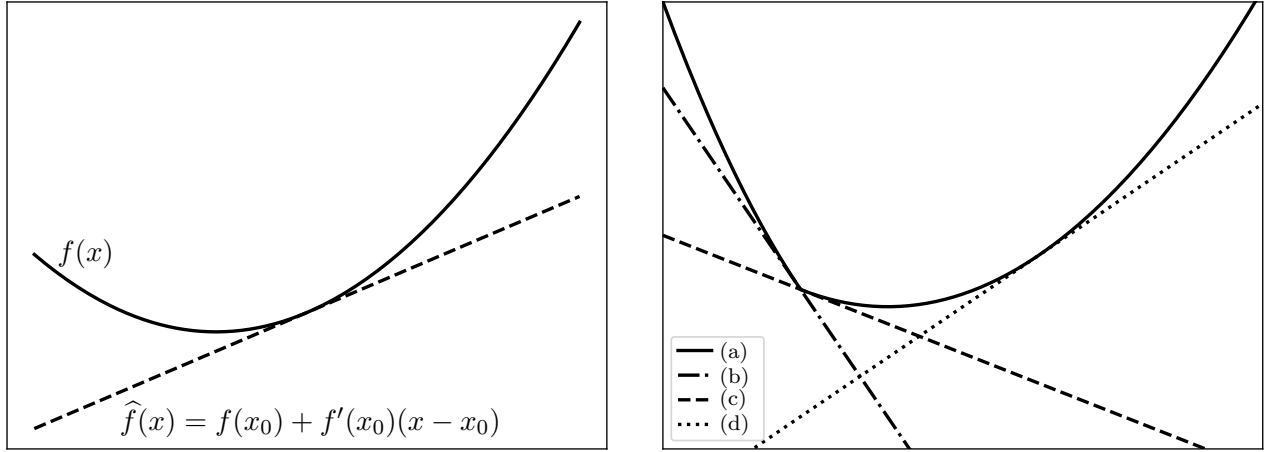


Figure 14.1. *Left:* The quadratic $f(x) = \frac{1}{2}x^2$ and the linear approximation $\hat{f}(x) = f(x_0) + s(x - x_0)$, where $x_0 = \frac{1}{2}$ and $s = f'(x_0)$. *Right:* the piecewise quadratic $f(x) = \max\{f_0(x), f_1(x)\}$ where $f_0(x) = \frac{1}{2}x^2$ and $f_1(x) = \frac{1}{4}(x + \frac{1}{4})^2 + \frac{1}{8}$, intersecting at $x_0 = \frac{1 - \sqrt{10}}{4}$. (a) The function $f(x)$. (b) The linear underestimator $\hat{f}(x) = f(x_0) + s_0(x - x_0)$ for $s_0 = f'_0(x_0)$. (c) The linear underestimator $\hat{f}(x) = f(x_0) + s_1(x - x_0)$ for $s_1 = f'_1(x_0)$. (d) The linear approximation $\hat{f}(x) = f(x_1) + f'(x_1)(x - x_1)$ around the point $x_1 = \frac{1}{4}$.

Lastly, we revisit a special duality relationship that all closed convex functions f enjoy (see Appendix C.2 for a fuller treatment). The *Fenchel-Legendre conjugate* or *convex conjugate* of a function f is

$$f^*(s) := \sup_x \{ \langle s, x \rangle - f(x) \}. \quad (14.1.3)$$

The function f^* is always convex, as it is the supremum of linear functions of s , and for any $x^*(s)$ maximizing $\langle s, x \rangle - f(x)$, we have that $x^*(s) \in \partial_s f^*(s)$ by the relationship (14.1.2); by a bit more work, we see that if $s \in \partial f(x)$, then $0 \in \partial_x \{ \langle s, x \rangle - f(x) \}$ and so x maximizes $\langle s, x \rangle - f(x)$. See Figure 14.2 for a graphical representation of this process. Flipping this argument by replacing f with f^* and x with s , when $s \in \partial f(x)$ and x maximizes $\langle s, x \rangle - f(x)$ in x , then $x \in \partial f^*(s)$ and so s maximizes $\langle s, x \rangle - f^*(s)$ in s . From this development comes the *biconjugate*, that is, $f^{**}(x) = \sup_s \{ \langle s, x \rangle - f^*(s) \}$, or $f^{**} = (f^*)^*$. The biconjugate f^{**} is equal to the supremum of *all* linear functionals below f :

Lemma 14.1.1. *Let f be a closed convex function. Then $f^{**}(x) = f(x)$.*

Proof We prove the equivalent statement that if $G \subset \mathbb{R}^d \times \mathbb{R}$ denotes all the pairs (s, b) so that the affine function $x \mapsto \langle s, x \rangle - b$ minorizes f , that is, $f(x) \geq \langle s, x \rangle - b$ for all x , then

$$f^{**}(x) = \sup_{(s,b) \in G} \{ \langle s, x \rangle - b \}.$$

Theorem B.3.7 in Appendix B.3.2 guarantees that $\sup_{(s,b) \in G} \{ \langle s, x \rangle - b \} = f(x)$. To see the displayed equality, note that $(s, b) \in G$ if and only if

$$f(x) \geq \langle s, x \rangle - b \text{ for all } x \text{ iff } b \geq \langle s, x \rangle - f(x) \text{ for all } x \text{ iff } b \geq f^*(s).$$

In particular, we obtain the equalities

$$\sup_{(s,b) \in G} \{ \langle s, x \rangle - b \} = \sup_{s,b} \{ \langle s, x \rangle - b \mid s \in \text{dom } f^*, f^*(s) \leq b \} = \sup_s \{ \langle s, x \rangle - f^*(s) \}$$

as desired. □

We immediately have the *Fenchel-Young inequality* that

$$f^*(s) + f(x) \geq \langle s, x \rangle \text{ for all } s, x,$$

and (see Proposition C.2.3) if f is a closed convex function, then equality holds if and only if

$$s \in \partial f(x) \text{ or } x \in \partial f^*(s), \quad (14.1.4)$$

which are equivalent. Thus we obtain the identities

$$\partial f^* = (\partial f)^{-1} \text{ and } \partial f = (\partial f^*)^{-1},$$

and we have the characterization

$$\partial f^*(s) = \operatorname{argmin}_x \{-\langle s, x \rangle + f(x)\} = \operatorname{argmax}_x \{\langle s, x \rangle - f(x)\}.$$

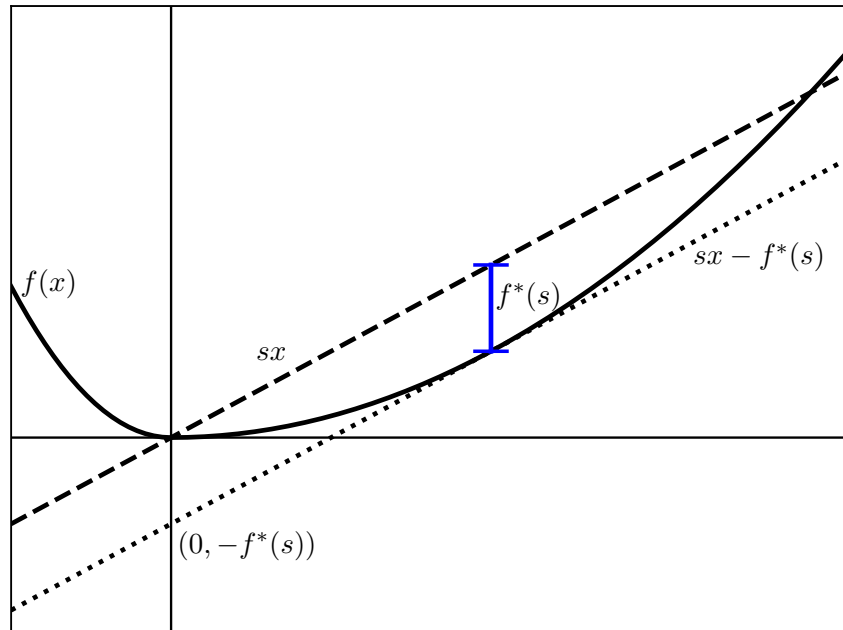


Figure 14.2. The conjugate function. The line of long dashes is $f(x) = sx$, while the dotted line is $x \mapsto sx - f^*(s)$. The blue line is the largest gap between sx and $f(x)$, which equals $f^*(s)$. Note that $x \mapsto sx - f^*(s)$ meets the graph of $f(x)$ at exactly the point of maximum difference $sx - f(x)$, where $f'(x) = s$.

14.1.2 From a proper loss to an entropy

The key construction underlying all of our proper losses is the optimal value of the expected loss. To any loss ℓ acting on a family \mathcal{P} of distributions, we construct the *generalized entropy* associated with the loss ℓ by

$$H_\ell(Y) := \inf_{Q \in \mathcal{P}} \mathbb{E}[\ell(Q, Y)], \quad (14.1.5)$$

where we have paralleled the typical notation $H(Y)$ for the Shannon entropy. In many cases, it will be more convenient to write this entropy directly as a function of the distribution P of Y , in which case we write

$$H_\ell(P) = \inf_{Q \in \mathcal{P}} \mathbb{E}_P[\ell(Q, Y)], \quad (14.1.6)$$

where Y follows the distribution P ; we will use whichever is more convenient. As the notation (14.1.6) makes clear, $H_\ell(P)$ is the infimum of a collection of linear functions of the form $P \mapsto \mathbb{E}_P[\ell(Q, Y)]$, one for each $Q \in \mathcal{P}$, so that necessarily $H_\ell(P)$ is concave in P . The remainder of this chapter, and several parts of the coming chapters, highlights the ways that this particular quantity informs the properties of the loss ℓ , and more generally, how we may always view any concave function H on a family of distributions \mathcal{P} as a generalized entropy function.

In Section 14.2, we show how such entropy-type functionals map back to losses themselves, so for now we content ourselves with a few examples to see why we call these entropies. Let us temporarily assume that Y has finite support $\{1, \dots, k\}$ with $\mathcal{P} = \Delta_k = \{p \in \mathbb{R}_+^k \mid \langle \mathbf{1}, p \rangle = 1\}$ the collection of probability mass functions on elements $\{1, \dots, k\}$.

Example 14.1.2 (Log loss): Consider the log loss $\ell_{\log}(p, y) = -\log p_y$. Then

$$H_{\ell_{\log}}(p) = \inf_{q \in \Delta_k} \mathbb{E}_p[-\log q_Y] = \inf_{q \in \Delta_k} \left\{ -\sum_{y=1}^k p_y \log \frac{q_y}{p_y} - \sum_{y=1}^k p_y \log p_y \right\} = -\sum_{y=1}^k p_y \log p_y,$$

the classical Shannon entropy. \diamond

This highlights an operational interpretation of entropy distinct from that arising in coding: the (Shannon) entropy is the minimal expected loss of a player in a prediction game, where the player chooses a distribution Q on Y , nature draws $Y \sim P$, and upon observing $Y = y$, the player suffers loss $-\log Q(Y = y)$.

Example 14.1.3 (0-1 error): If instead we take the 0-1 loss, that is, $\ell_{0-1}(p, y) = 1$ if $p_y \leq p_j$ for some $j \neq y$ and $\ell_{0-1}(p, y) = 0$ otherwise, then

$$H_{\ell_{0-1}}(p) = \inf_{q \in \Delta_k} \mathbb{E}_p[\ell(q, y)] = 1 - \max_y p_y.$$

So $H_{\ell_{0-1}}(e_y) = 0$ for any standard basis vector, that is, distribution with all mass on a single point y , and $H_{\ell_{0-1}}(p) > 0$ otherwise. Moreover, the vector $p = \mathbf{1}/k$ maximizes $H_{\ell_{0-1}}(p)$, with $H_{\ell_{0-1}}(\mathbf{1}/k) = \frac{k-1}{k}$. \diamond

Example 14.1.4 (Brier scoring and squared error): For the squared error (Brier scoring) loss $\ell_{\text{sq}}(p, y) = \|p - e_y\|_2^2$, where $e_y \in \{0, 1\}^k$ is the y th standard basis vector, let Y have p.m.f. $p \in \Delta_k$. Then

$$H_{\ell_{\text{sq}}}(Y) = \mathbb{E}[\ell_{\text{sq}}(p, Y)] = \|p\|_2^2 - 2\|p\|_2^2 + 1 = 1 - \|p\|_2^2.$$

So as above, we have $H_{\ell_{\text{sq}}}(Y) \geq 0$, with $H_{\ell_{\text{sq}}}(Y) = 0$ if and only if Y is a point mass on one of $\{1, \dots, k\}$, and the uniform distribution with p.m.f. $p = \frac{1}{k}\mathbf{1}$ maximizes the entropy, with $H_{\ell_{\text{sq}}}(\text{Uniform}([k])) = 1 - 1/k$. \diamond

These examples highlight how these entropy functions are types of uncertainty measures, giving rise to “maximally uncertain” distributions p , which are typically uniform on Y .

14.1.3 The information in an experiment

In classical information theory, the mutual (Shannon) information between random variables X and Y is the gap between the entropy of Y and the remaining entropy given X , that is,

$$I(X; Y) = H(Y) - H(Y | X).$$

In complete analogy with our development in Chapter 2, then, we can define the information between variables X and Y relative to a particular loss function ℓ . Thus, we define the ℓ -conditional entropy

$$H_\ell(Y | X = x) := \inf_{Q \in \mathcal{P}} \mathbb{E}[\ell(Q, Y) | X = x]$$

and, in analogy to the definitions in Section 2.1.1, the conditional entropy of Y given X is

$$H_\ell(Y | X) := \mathbb{E} \left[\inf_{Q \in \mathcal{P}} \mathbb{E}[\ell(Q, Y) | X] \right] = \int_{\mathcal{X}} H_\ell(Y | X = x) dP(x),$$

the average minimal expected loss when one observes X .

With this definition, we then can discuss the *information in an experiment*. This nomenclature follows classical statistical parlance, where by an experiment, we mean the observation of a variable X in a Markov chain $X \rightarrow Y$, where we think of Y as a hypothesis to be tested or a value to be predicted, and we ask how much observing X helps to actually allow this prediction. Then we define

$$I_\ell(X; Y) := H_\ell(Y) - H_\ell(Y | X), \quad (14.1.7)$$

which is nonnegative and is the gap between the prior entropy of Y and its posterior entropy conditional on the observation X . That is, this information measure is precisely the gap between the best achievable loss in the prediction of a distribution P for Y *a priori*, when we observe nothing, and that achievable *a posteriori*, when we observe X . In parallel to our alternative view of the entropy as the (expected) minimal loss of a player in a prediction game, then, the information between X and Y is the improvement an observation X offers a player in predicting Y when measuring error with the loss ℓ . The information (14.1.7) is typically asymmetrical in X and Y , so we are careful about the ordering (this lack of symmetric holds, essentially, unless ℓ is the log loss).

The next three examples show different information quantities, where in each we let \mathcal{Y} have finite cardinality k , and thus identify \mathcal{P} with the probability simplex $\Delta_k = \{p \in \mathbb{R}_+^k \mid \langle \mathbf{1}, p \rangle = 1\}$.

Example 14.1.5 (Shannon information): Taking the log loss $\ell(p, y) = -\log p_y$, we have

$$I_\ell(X; Y) = H_\ell(Y) - H_\ell(Y | X) = H(Y) - H(Y | X) = I(X; Y),$$

the classical Shannon information. \diamond

Example 14.1.6 (0-1 error): Consider the 0-1 error $\ell_{0-1}(p, y) = 1$ if $p_y \leq \max_{j \neq y} p_j$ and $\ell_{0-1}(p, y) = 0$ if $p_y > \max_{j \neq y} p_j$. Then letting $y^* = \operatorname{argmax}_y \mathbb{P}(Y = y)$ and $y^*(x) = \operatorname{argmax}_y \mathbb{P}(Y = y | X = x)$, we have

$$I_{\ell_{0-1}}(X; Y) = \mathbb{P}(Y = y^*) - \mathbb{E}[\mathbb{P}(Y = y^*(X) | X)] = \mathbb{P}(Y = y^*) - \mathbb{P}(Y = y^*(X)),$$

the gap between the prior probability of making a mistake when guessing Y and the posterior probability given X . \diamond

Example 14.1.7 (Squared error): For the Brier score with squared error $\ell_{\text{sq}}(p, y) = \|p - e_y\|_2^2$, we have $H_{\ell_{\text{sq}}}(p) = 1 - \|p\|_2^2$, and so

$$I_{\ell_{\text{sq}}}(X; Y) = \sum_{j=1}^k \mathbb{E} [\mathbb{P}(Y = j | X)^2] - \sum_{j=1}^k \mathbb{P}(Y = j)^2 = \sum_{j=1}^k \text{Var}(\mathbb{P}(Y = j | X)),$$

the summed variances of the random variables $\mathbb{P}(Y = j | X)$. The higher the variance of these quantities, the more information X carries about Y . \diamond

14.2 Characterizing proper losses and Bregman divergences

With the definition (14.1.5) of the fundamental generalized entropy, we can now proceed to a characterization of all proper losses. We do this in three settings: in the first (Section 14.2.1), we give a representation for proper losses when \mathcal{Y} is finite and discrete, so we can identify it with $\mathcal{Y} = \{1, \dots, k\}$ and distributions P on \mathcal{Y} with probability mass functions $p \in \Delta_k$. We then demonstrate a full characterization of propriety (Section 14.2.2), which requires measure-theoretic tools and can be skipped. As the final approach to considering propriety, we modify the results for finite \mathcal{Y} to consider cases in which Y is vector-valued and $\mathcal{Y} \subset \mathbb{R}^k$ is contained in a compact set. This case transparently generalizes the finite representations of Section 14.2.1 and will form the basis of our development going forward, as it allows us to more directly apply to tools of convexity and analysis.

14.2.1 Characterizing proper losses for Y taking finitely many values

Here, we present the *Savage representation* of proper losses, which characterizes all proper losses using the entropies (14.1.5) or, equivalently, (14.1.6). To avoid pathological cases, we work with regular losses, which always assign a finite value to the correct predicted distribution; we assume regularity without further comment.

Definition 14.2. Let \mathcal{P} be a family of distribution on \mathcal{Y} . The loss $\ell : \mathcal{P} \times \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ is regular for the family \mathcal{P} if $\mathbb{E}_P[\ell(P, Y)]$ is real valued for all $P \in \mathcal{P}$.

We do allow losses to attain infinite values, for example, we can allow $\ell(Q, y) = +\infty$ if Q assigns probability 0 to an event y , as in the case of the logarithmic loss. The following theorem then provides the promised representation of proper losses, and additionally, highlights the centrality of the generalized entropy functionals.

Theorem 14.2.1 (Proper scoring rules: the finite case). Let $\mathcal{Y} = \{1, \dots, k\}$ be finite and $\mathcal{P} \subset \Delta_k$ a convex collection of distributions on \mathcal{Y} . Then the following are true.

(i) If the loss $\ell : \mathcal{P} \times \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ satisfies the representation

$$\ell(p, y) = -\Omega(p) - \langle \nabla \Omega(p), e_y - p \rangle \quad (14.2.1)$$

for a subdifferentiable closed convex function $\Omega : \mathcal{P} \rightarrow \mathbb{R}$, where $\nabla \Omega(p) \in \partial \Omega(p)$, then ℓ is proper.

(ii) Conversely, if ℓ is proper, then choosing Ω to be the negative generalized entropy

$$\Omega_\ell(p) := -H_\ell(p) = \sup_q \{-\mathbb{E}_p[\ell(q, Y)] \mid q \in \mathcal{P}\}$$

satisfies equality (14.2.1) (and h is closed).

Additionally, if ℓ is real valued, then $\nabla\Omega(p) \in \mathbb{R}^k$ in the representation (14.2.1). If $\ell(p, y)$ can take the value $+\infty$, then we allow $\nabla\Omega(p) \in \mathbb{R}^k$ when $p \notin \text{relint } \Delta_k$. The loss is strictly proper if and only if the convex Ω is strictly convex.

Proof If ℓ has the given representation and $\mathbb{P}(Y = y) = p_y$, then we have

$$\mathbb{E}[\ell(q, Y)] = -\Omega(q) - \langle \nabla\Omega(q), p - q \rangle \geq -\Omega(p) = \mathbb{E}[\ell(p, Y)]$$

by the first-order convexity property of convex functions (that is, the definition (14.1.1)) of a subdifferential).

Conversely, suppose that the loss is proper, and let $\Omega(p) = \Omega_\ell(p)$. Clearly Ω is convex, as it is the supremum of linear functionals of p . Moreover, propriety of ℓ guarantees that

$$\Omega(p) \geq -\mathbb{E}[\ell(q, Y)] = \Omega(q) + \sum_{y=1}^k -\ell(q, y)(p_k - q_k)$$

That is, for each $q \in \mathcal{P}$ the vector $[-\ell(q, y)]_{y=1}^k \in \partial\Omega(q)$, so Ω is subdifferentiable. Choosing the vector $\nabla\Omega(p) = [-\ell(p, y)]_{y=1}^k$, we have

$$\ell(p, y) = -\Omega(p) + \ell(p, y) + \Omega(p) = -\Omega(p) - \sum_{i=1}^k p_i \ell(p, i) + \ell(p, y) = -\Omega(p) - \langle \nabla\Omega(p), e_y - p \rangle$$

as desired. Note that $\ell(p, y) < \infty$ except when $p_y = 0$, in which case our definition $\nabla\Omega(p) = [-\ell(p, y)]_{y=1}^k$ remains sensible as $-\langle \nabla\Omega(p), e_y - p \rangle = +\infty$.

As an alternative argument more directly using convexity, definition of $\Omega(p) = \sup_q \{-\mathbb{E}_p[\ell(q, Y)] \mid q \in \mathcal{P}\}$ and the immediate calculation (14.1.2) of the subdifferential of the supremum shows that

$$\partial\Omega(p) \supset \left\{ [-\ell(q, y)]_{y=1}^k \mid q \in \Delta_k \text{ satisfies } -\mathbb{E}_p[\ell(q, Y)] = \Omega(p) \right\}.$$

But propriety guarantees that the set of such q includes p , so that $\partial\Omega(p) \supset [-\ell(p, y)]_{y=1}^k$.

For the strict inequalities and strict propriety, trace the argument replacing inequalities with strict inequalities for $q \neq p$ and use Corollary B.3.2 or C.1.9. \square

The negative generalized entropy Ω in Theorem 14.2.1 is essentially unique and marks an important duality between proper losses and convex functions: to each loss, we can assign a generalized entropy, and from this generalized entropy, we can reconstruct the loss. Exercise 14.2 explores this connection. We can also give a few examples that show how to recover standard losses. For each, we begin with a convex function Ω , then exhibit the associated proper or strictly proper scoring rule. One thing to notice in this representation is that, typically, we do *not* expect to achieve a loss function convex in p , which is a weakness of the representation (14.2.1). In Section 14.3 (and Chapter 16, especially section 16.5) in more depth), however, we will show how to convert suitable proper losses into *surrogates* that are convex in their arguments and which, after a particular transformation based on convex duality, are proper and yield the correct distributional predictions. We defer this, however, and instead provide a few examples.

Example 14.2.2 (Logarithmic losses): Consider the negative entropy $\Omega(p) = \sum_{y=1}^k p_y \log p_y$. We have $\frac{\partial}{\partial p_y} \Omega(p) = 1 + \log p_y \in [-\infty, 1]$, and

$$\ell_{\log}(p, y) = - \sum_{j=1}^k p_j \log p_j + \sum_{j=1}^k p_y (1 + \log p_j) - (1 + \log p_y) = -\log p_y,$$

yielding the log loss. Note that for this case, we do require that the gradients $\nabla \Omega(p)$ take values in the (downward) extended reals $\underline{\mathbb{R}}^k$. \diamond

Example 14.2.3 (Brier scores and squared error): When we have the squared error $\ell_{\text{sq}}(p, y) = \|p - e_y\|_2^2$, we can directly check that $\Omega(p) = \|p\|_2^2$ gives the loss. Indeed,

$$-\|p\|_2^2 - 2\langle p, e_y - p \rangle = \|p\|_2^2 - 2\langle p, e_y \rangle + 1 - 1 = \|p - e_y\|_2^2 - 1.$$

So aside from an additive constant, we have the desired result. \diamond

More esoteric examples exist in the literature, such as the spherical score arising from $\Omega(p) = \|p\|_2$ (note the lack of a square).

Example 14.2.4 (Spherical scores): Let $\Omega(p) = \|p\|_2$, which is strictly convex on Δ_k . Then

$$\nabla \Omega(p) = p / \|p\|_2$$

and $\ell(p, y) = -\|p\|_2 - \frac{1}{\|p\|_2} \langle p, e_y - p \rangle = -p_y / \|p\|_2$, which is strictly proper but does not retain convexity. \diamond

Bregman divergences

A key aspect of the Savage representation (14.2.1) is that associated to any proper loss is a *Bregman divergence*, which measures the difference between a convex function and its first-order approximation. Recall from Chapter 3 that for a function $\Omega : \mathbb{R}^k \rightarrow \overline{\mathbb{R}}$, we define

$$D_\Omega(u, v) := \Omega(u) - \Omega(v) - \langle \nabla \Omega(v), u - v \rangle. \quad (14.2.2)$$

In typical definitions of the divergence, one requires that Ω be differentiable; here, we allow non-differentiable Ω so long as the choice $\nabla \Omega(v) \in \partial \Omega(v)$ is given. In particular, we see that

$$D_\Omega(u, v) \geq 0$$

for all u and v , and moreover, if Ω is strictly convex

$$D_\Omega(u, v) > 0 \text{ whenever } u \neq v.$$

(See, e.g., Corollaries B.3.2 and C.1.9 in the appendices.)

Familiar examples include the squared Euclidean norm $\Omega(u) = \frac{1}{2} \|u\|_2^2$, which by inspection gives

$$D_\Omega(u, v) = \frac{1}{2} \|u - v\|_2^2,$$

and the negative entropies $\Omega(u) = \sum_{j=1}^k u_j \log u_j$, which implicitly encodes the constraint that $u \succ 0$. This gives

$$D_\Omega(u, v) = \sum_{j=1}^k u_j \log u_j - \sum_{j=1}^k v_j \log v_j - \sum_{j=1}^k (1 + \log v_j)(u_j - v_j) = \sum_{j=1}^k u_j \log \frac{u_j}{v_j} + \mathbf{1}^T(u - v).$$

If $u, v \in \Delta_k$, then evidently $D_\Omega(u, v) = D_{\text{kl}}(u \| v)$ because $\mathbf{1}^T u = \mathbf{1}^T v = 1$, where we identify u and v with probability mass functions.

Continuing this identification of distributions on \mathcal{Y} with elements $p \in \Delta_k$ in the probability simplex, we can reconsider the gaps between a loss evaluated at a true distribution p and an alternative q . In this case, the representation Theorem 14.2.1 provides allows us to connect proper losses with first-order divergences immediately. Indeed, let $\Omega : \Delta_k \rightarrow \bar{\mathbb{R}}$ be a convex function and loss ℓ be the associated proper loss, with $\ell(p, y) = -\Omega(p) - \langle \nabla \Omega(p), e_y - p \rangle$. Now, suppose that Y has p.m.f. p ; then for any $q \in \Delta_k$, the gap

$$\begin{aligned} \mathbb{E}_p[\ell(q, Y)] - \mathbb{E}_p[\ell(p, Y)] &= \Omega(p) - \Omega(q) - \sum_{y=1}^k p_y \langle \nabla \Omega(q), e_y - q \rangle \\ &= \Omega(p) - \Omega(q) - \langle \nabla \Omega(q), p - q \rangle = D_\Omega(p, q). \end{aligned}$$

We record this as a corollary to Theorem 14.2.1, highlighting the links between propriety, first-order divergences, and proper loss functions.

Corollary 14.2.5. *Let the conditions of Theorem 14.2.1 hold. Then ℓ is (strictly) proper if and only if there exists a (strictly) convex $\Omega : \Delta_k \rightarrow \bar{\mathbb{R}}$ for which*

$$\mathbb{E}_p[\ell(q, Y)] - \mathbb{E}_p[\ell(p, Y)] = D_\Omega(p, q)$$

for all $p, q \in \Delta_k$.

14.2.2 General proper losses

More generally, we can consider predicting distributions P on general sets \mathcal{Y} . For example, recalling the meteorological motivation of predicting the weather, suppose we wish to predict a distribution of the (real-valued) amount Y of rainfall on a given day. Many predictions place a point mass at $Y = 0$, with a decaying tail for higher amounts of rainfall. Then it is natural to predict a cumulative distribution function $F : \mathbb{R} \rightarrow [0, 1]$, measuring error relative to the actual amount of rain that falls. Several losses are common in the literature; one common example is the *continuous ranked probability score*.

Example 14.2.6 (Continuous ranked probability score (CRPS)): The CRPS loss for a CDF F at y is

$$\ell_{\text{crps}}(F, y) = \int (F(t) - \mathbf{1}\{y \leq t\})^2 dt. \quad (14.2.3)$$

This is a strictly proper scoring rule: let G be any cumulative distribution function, meaning that $\lim_{t \rightarrow -\infty} G(t) = 0$ and $\lim_{t \rightarrow \infty} G(t) = 1$, and let Y have CDF F . Then

$$\begin{aligned} \mathbb{E}[\ell_{\text{crps}}(G, Y)] - \mathbb{E}[\ell_{\text{crps}}(F, Y)] &= \int (G(t)^2 - F(t)^2 - 2(G(t) - F(t))\mathbb{E}[\mathbf{1}\{Y \leq t\}]) dt \\ &= \int (G(t) - F(t))^2 dt \end{aligned}$$

because $\mathbb{E}[\mathbf{1}\{Y \leq t\}] = F(t)$. This is a variant of the (squared) Cramér-von-Mises distance between F and G , and which is positive unless $F = G$. Unfortunately, computing the CRPS loss (14.2.3) is often challenging except for specially structured F . \diamond

Because the computation of the continuous ranked probability score is challenging, it can be advantageous to consider other losses on probability distributions, which can allow more flexibility in modeling. To that end, we define the *quantile loss*: for a probability distribution P on Y , let

$$\text{Quant}_\alpha(Y) = \text{Quant}_\alpha(P) := \inf \{t \mid P(Y \leq t) \geq \alpha\}$$

to be the α -quantile of the distribution P . (When Y has cumulative distribution F , this is the inverse CDF mapping $F^{-1}(\alpha) = \inf\{t \mid F(t) \geq \alpha\}$.) Defining the quantile penalty

$$\rho_\alpha(t) = \alpha [t]_+ + (1 - \alpha) [-t]_+,$$

for a collection \mathcal{A} of values in $[0, 1]$, the *quantile loss* is

$$\ell_{\text{quant}, \mathcal{A}}(P, y) := \sum_{\alpha \in \mathcal{A}} \rho_\alpha(y - \text{Quant}_\alpha(P)). \quad (14.2.4)$$

The propriety of the quantile loss is relatively straightforward; it is, however, not strictly proper.

Example 14.2.7 (Quantile loss): To see that the quantile loss (14.2.4) is proper, consider the single quantile penalty ρ_α : let $g(t) = \mathbb{E}[\rho_\alpha(Y - t)] = \alpha \mathbb{E}[[Y - t]_+] + (1 - \alpha) \mathbb{E}[[t - Y]_+]$, which we claim is minimized by $\text{Quant}_\alpha(Y)$. Indeed, g is convex, and it has left and right derivatives

$$\begin{aligned} \partial_- g(t) &:= \lim_{s \uparrow t} \frac{g(s) - g(t)}{s - t} = -\alpha \mathbb{P}(Y \geq t) + (1 - \alpha) \mathbb{P}(Y < t) = \mathbb{P}(Y < t) - \alpha \quad \text{and} \\ \partial_+ g(t) &:= \lim_{s \downarrow t} \frac{g(s) - g(t)}{s - t} = -\alpha \mathbb{P}(Y > t) + (1 - \alpha) \mathbb{P}(Y \leq t) = \mathbb{P}(Y \leq t) - \alpha. \end{aligned}$$

Indeed, for $t = \text{Quant}_\alpha(Y)$, we have $\partial_- g(t) = \mathbb{P}(Y < t) - \alpha \leq 0$ and $\partial_+ g(t) = \mathbb{P}(Y \leq t) - \alpha \geq 0$, because $t \mapsto \mathbb{P}(Y \leq t)$ is right continuous. So convexity yields

$$\mathbb{E}[\rho_\alpha(Y - \text{Quant}_\alpha(Y))] \leq \mathbb{E}[\rho_\alpha(Y - t)]$$

for all t . Applying this argument for each $\alpha \in \mathcal{A}$, we thus have

$$\mathbb{E}[\ell_{\text{quant}, \mathcal{A}}(Q, Y)] \geq \mathbb{E}[\ell_{\text{quant}, \mathcal{A}}(P, Y)]$$

for any Q whenever $Y \sim P$, and equality holds whenever Q and P have identical α quantile for each $\alpha \in \mathcal{A}$. \diamond

The general case of Theorem 14.2.1 allows us to address such scenarios, though it does require measure theory to properly define. Happily, the generality does not require a particularly more sophisticated proof. For a (convex) function $\Omega : \mathcal{P} \rightarrow \overline{\mathbb{R}}$ on a family of distributions \mathcal{P} on a set \mathcal{Y} , we say $\Omega'(P; \cdot) : \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ is a subderivative of Ω at $P \in \mathcal{P}$ whenever

$$\begin{aligned} \Omega(Q) &\geq \Omega(P) + \int_{\mathcal{Y}} \Omega'(P, y)(dQ(y) - dP(y)) \quad \text{for all } Q \in \mathcal{P}. \\ &= \Omega(P) + \mathbb{E}_Q[\Omega'(P, Y)] - \mathbb{E}_P[\Omega'(P, Y)] \end{aligned} \quad (14.2.5)$$

When \mathcal{Y} is discrete and we can identify \mathcal{P} with the simplex Δ_k , the inequality (14.2.5) is simply the typical subgradient inequality (14.1.1) that $\Omega(q) \geq \Omega(p) + \langle \nabla \Omega(p), q - p \rangle$ for $p, q \in \Delta_k$, where $\nabla \Omega(p) \in \partial \Omega(p)$. We then have the following generalization of Theorem 14.2.1.

Theorem 14.2.8. *Let \mathcal{P} be a convex collection of distributions on \mathcal{Y} . Then the following are true.*

(i) *If the loss $\ell : \mathcal{P} \times \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ satisfies the representation*

$$\ell(P, y_0) = -\Omega(P) - \Omega'(P, y_0) + \int \Omega'(P, y) dP(y), \quad \text{for all } y_0 \in \mathcal{Y}, \quad (14.2.6)$$

where $\Omega'(P, \cdot) : \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ is a subderivative of Ω at $P \in \mathcal{P}$, then it is proper.

(ii) *Conversely, if ℓ is proper, then choosing Ω to be the negative generalized entropy $\Omega_\ell(P) = -H_\ell(P) = \sup\{-\mathbb{E}_P[\ell(Q, Y)] \mid Q \in \mathcal{P}\}$ satisfies equality (14.2.6).*

The loss is strictly proper if and only if the convex Ω is strictly convex.

Proof If ℓ has the representation (14.2.6), then we have

$$-\mathbb{E}_P[\ell(P, Y)] = \Omega(P) \geq \Omega(Q) + \int \Omega'(Q, y)(dP(y) - dQ(y)) = -\mathbb{E}_P[\ell(Q, Y)]$$

for any $Q \in \mathcal{P}$ by the definition (14.2.5) of a subderivative. Rewriting, we have $\mathbb{E}_P[\ell(P, Y)] \leq \mathbb{E}_P[\ell(Q, Y)]$ and ℓ is proper.

Conversely, if ℓ is proper and regular, then as in the proof of Theorem 14.2.1 we define

$$\Omega(P) := \sup_{Q \in \mathcal{P}} -\mathbb{E}_P[\ell(Q, Y)] = -\mathbb{E}_P[\ell(P, Y)],$$

which is the supremum of linear functionals of P and hence convex. If we let $\Omega'(P, y) = -\ell(P, y) \in \overline{\mathbb{R}}$ for $P \in \mathcal{P}$, then

$$\Omega(P) \geq -\mathbb{E}_P[\ell(Q, Y)] = \Omega(Q) + \mathbb{E}_Q[\ell(Q, Y)] - \mathbb{E}_P[\ell(Q, Y)] = \Omega(Q) + \int \Omega'(P, y)(dP(y) - dQ(y))$$

by propriety, so that evidently $\Omega'(P, y)$ is a subderivative of Ω at $P \in \mathcal{P}$. That $L(P, y_0) = -\Omega(P) - \Omega'(P, y_0) + \int \Omega'(P, y) dP(y)$ is then immediate.

The arguments for strict propriety/convexity are similar. \square

The obvious corollary to Theorem 14.2.8, paralleling Corollary 14.2.5, follows.

Corollary 14.2.9. *Let \mathcal{P} be a convex collection of probability distributions on \mathcal{Y} . Then the loss $\ell : \mathcal{P} \times \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ is proper if and only if there exists a convex function $\Omega : \mathcal{P} \rightarrow \overline{\mathbb{R}}$ with subderivatives $\Omega'(P, \cdot) : \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ such that*

$$\ell(P, y_0) = -\Omega(P) - \Omega'(P, y_0) + \mathbb{E}_P[\Omega'(P, Y)] \quad \text{for all } y_0 \in \mathcal{Y}.$$

The loss ℓ is strictly proper if and only if Ω is strictly concave. Similarly, ℓ is (strictly) proper if and only if there exists a (strictly) convex and subdifferentiable $\Omega : \mathcal{P} \rightarrow \overline{\mathbb{R}}$ for which

$$\mathbb{E}_P[\ell(Q, Y)] - \mathbb{E}_P[\ell(P, Y)] = D_\Omega(P, Q).$$

The subdifferentials and differentiability in this potentially infinite dimensional case can make writing the particular representation (14.2.6) challenging; for example, the representation of the quantile loss in Example 14.2.7 is quite complex. In the case of predictions involving the cumulative distribution function F , however, one can obtain the subderivative by taking directional (Gateaux) derivatives in directions $G - F$ for cumulative distributions G . In this case, for the point cumulative distribution G_y with $G_y(t) = \mathbf{1}\{y \leq t\}$, we define

$$\Omega'(F, y) = \lim_{\epsilon \downarrow 0} \frac{\Omega(F + \epsilon(G_y - F)) - \Omega(F)}{\epsilon}.$$

The continuous ranked probability score (Example 14.2.6) admits this expansion.

Example 14.2.10 (CRPS (Example 14.2.6) continued): The strict propriety of the CRPS loss (14.2.3) means that the generalized negative entropy

$$\Omega(F) = \sup_G -\mathbb{E}[\ell(G, Y)] = -\mathbb{E}[\ell_{\text{crps}}(F, Y)] = \int (F(t) - 1)F(t)dt$$

by definition. Expanding $\Omega(F + \epsilon(G - F))$ for small ϵ as in the recipe above, we have

$$\Omega(F + \epsilon(G - F)) = \Omega(F) - \epsilon \int (G(t) - F(t))dt + 2\epsilon \int F(t)(G(t) - F(t))dt + O(\epsilon^2).$$

to obtain the y -based derivative $\Omega'(F, y)$, we choose $G_y(t) = \mathbf{1}\{y \leq t\}$ to obtain directional derivative

$$\Omega'(F, y) = \lim_{\epsilon \downarrow 0} \frac{\Omega(F + \epsilon(G_y - F)) - \Omega(F)}{\epsilon} = \int (\mathbf{1}\{y \leq t\} - F(t))dt - 2 \int F(t)(\mathbf{1}\{y \leq t\} - F(t))dt.$$

By inspection, when Y has cumulative distribution function F , $\mathbb{E}[\Omega'(F, Y)] = 0$ and so

$$\begin{aligned} & -\Omega(F) - \Omega'(F, y) + \mathbb{E}[\Omega'(F, Y)] \\ &= \int (-F(t)^2 + F(t) - F(t) + \mathbf{1}\{y \leq t\} + 2F(t)^2 - 2F(t)\mathbf{1}\{y \leq t\}) dt \\ &= - \int (F(t) - \mathbf{1}\{y \leq t\})^2 dt = \ell_{\text{crps}}(F, y), \end{aligned}$$

as desired. \diamond

We can also consider the log loss for general distributions, which is a bit subtle because of the necessity of defining a base measure. Note that for discrete cases—when $Y \in \{1, \dots, k\}$ or \mathcal{Y} countable—we could always use the probability mass function without loss of generality.

Example 14.2.11 (Log loss for general distributions): For general distributions P , the logarithmic loss requires additional work to be sensible, because $-\log p(y)$ is not well-defined if \mathcal{X} is not discrete. Thus, we fix a base measure ν on \mathcal{Y} , and let \mathcal{P} be the collection of distributions that are absolutely continuous with respect to ν . Then for $P \in \mathcal{P}$, we let $\ell(P, y) = -\log p(y)$, where p is the density of P with respect to ν , while $\ell(P, y) = +\infty$ for $P \not\ll \nu$. Noting that for $P \in \mathcal{P}$,

$$\mathbb{E}_P[\ell(Q, Y)] = \mathbb{E}_P[-\log q(Y)] = \mathbb{E}_P \left[\log \frac{p(Y)}{q(Y)} \right] - \mathbb{E}_P[\log p(Y)] = D_{\text{kl}}(P \| Q) + H_\nu(P),$$

where $H_\nu(P) := -\int p(y) \log p(y) d\nu(y)$ denotes the (Shannon) entropy for the base measure ν , we see that ℓ is indeed strictly proper.

The negative entropy is thus $\Omega(P) = -H_\nu(P)$ when $P \ll \nu$ and $+\infty$ otherwise. To obtain the directional derivatives as in the representation (14.2.6), we heuristically take the derivative $\frac{\partial}{\partial p(y)} \Omega(P)$ to guess that

$$\Omega'(P, y) = 1 + \log p(y).$$

We can directly check that this is indeed a subgradient for Ω : we have

$$\begin{aligned} \Omega(P) + \mathbb{E}_Q[\Omega'(P, Y)] - \mathbb{E}_P[\Omega'(P, Y)] &= \mathbb{E}_P[\log p(Y)] + \mathbb{E}_Q[1 + \log p(Y)] - \mathbb{E}_P[1 + \log p(Y)] \\ &= \mathbb{E}_Q \left[\log \frac{p(Y)}{q(Y)} \right] + \mathbb{E}_Q[\log q(Y)] = \Omega(Q) - D_{\text{kl}}(Q \| P), \end{aligned}$$

assuming $Q \ll \nu$. (Otherwise, $D_{\text{kl}}(Q \| P) = +\infty$, and $\mathbb{E}_Q[\log p(Y)] = -\infty$ regardless.) Finally, we check the loss representation (14.2.6): so long as $Q \ll \nu$, we have

$$-\Omega(Q) - \Omega'(Q, y) + \mathbb{E}_Q[\Omega'(Q, Y)] = H_\nu(Q) - H_\nu(Q) - 1 + 1 - \log q(y) = -\log q(y),$$

as desired. \diamond

14.2.3 Proper losses and vector-valued Y

The final variant of propriety we consider generalizes that when \mathcal{Y} is finite and identified with $\{1, \dots, k\}$ in Section 14.2.1. Now, we assume that Y is vector-valued, with $\mathcal{Y} \subset \mathbb{R}^k$, and assume the convex hull

$$\text{Conv}(\mathcal{Y}) = \{\mathbb{E}_P[Y] \mid P \text{ is a distribution on } \mathcal{Y}\}$$

is bounded. (Typically, it will also be compact, though this will not be central to our development, and pathological cases, such as $\mathcal{Y} = \{1/n\}_{n \in \mathbb{N}}$, exist.) An example showing how to use this representation for multinomial $Y \in \{1, \dots, k\}$ may be clarifying.

Example 14.2.12 (Multinomial Y as vectors): If Y is a multinomial taking values in a discrete set of size k , we can instead identify Y with the first k standard basis vectors e_1, \dots, e_k . Then $p = \mathbb{E}[Y] \in \Delta_k$ is the p.m.f. of Y , and $\text{Conv}(\mathcal{Y}) = \Delta_k$. \diamond

Example 14.2.13 (Binary Y as a scalar): When $Y \in \{0, 1\}$ is a Bernoulli random variable, we identify Y with itself, so that $p = \mathbb{E}[Y] = \mathbb{P}(Y = 1) \in [0, 1]$ and $\text{Conv}(\mathcal{Y}) = [0, 1]$. \diamond

Example 14.2.14 (Ordinal Y as a scalar): Consider a rating problem of predicting the rating Y of a movie from 1 to 5 stars. In this case, Y takes values $\{1, \dots, 5\} \subset \mathbb{R}$, but the ordering between the elements is important; it is unnatural to treat Y as a multinomial. More generally, Y may take values in $\{y_1, \dots, y_k\} \subset \mathbb{R}$, where $y_1 < \dots < y_k$. As in the binary case, we identify Y with its scalar value, so that $\mathbb{E}[Y] \in [y_1, y_k]$ and $\text{Conv}(\mathcal{Y}) = [y_1, y_k]$. \diamond

In this vector-valued Y case, instead of prediction distributions P , the goal is to predict the *mean mapping*

$$\mu(P) := \mathbb{E}_P[Y] \in \text{Conv}(\mathcal{Y}),$$

so that $\mu : \mathcal{P} \rightarrow \mathbb{R}^k$ for the collection \mathcal{P} of distributions on Y . Our goal is to reward predictions of the correct expectation, leading to the following definition.

Definition 14.3. Let $\mathcal{M} = \text{cl Conv}(\mathcal{Y})$ be the (convex) collection of mean parameters for $Y \in \mathcal{Y}$. Then $\ell : \mathcal{M} \times \mathcal{Y} \rightarrow \bar{\mathbb{R}}$ is proper if

$$\mathbb{E}_P[\ell(\mu, Y)] \geq \mathbb{E}_P[\ell(\mathbb{E}_P[Y], Y)] \text{ for all } \mu \in \mathcal{M},$$

and strictly proper if the inequality is strict whenever $\mu \neq \mathbb{E}_P[Y]$.

Definition 14.3 generalizes Definition 14.2 in the multinomial case, where \mathcal{Y} is a discrete set that we may identify with the basis vectors $\{e_1, \dots, e_k\}$, as Example 14.2.12 makes clear.

With this definition, we can extend Theorem 14.2.1 to a more general case, where as usual we say that ℓ is regular if $\mathbb{E}_P[\ell(\mathbb{E}_P[Y], Y)] < \infty$ for all distributions P on \mathcal{Y} .

Theorem 14.2.15. Let $\mathcal{Y} \subset \mathbb{R}^k$ be finite, \mathcal{P} be the collection of distributions on \mathcal{Y} , and $\mathcal{M} = \text{Conv}(\mathcal{Y}) = \{\mathbb{E}_P[Y] \mid P \in \mathcal{P}\}$. A regular loss $\ell : \mathcal{M} \times \mathcal{Y} \rightarrow \bar{\mathbb{R}}$ is proper if and only if there exists a closed convex $\Omega : \mathcal{M} \rightarrow \bar{\mathbb{R}}$ such that

$$\ell(\mu, y) = -\Omega(\mu) - \langle \nabla \Omega(\mu), y - \mu \rangle$$

for some subgradient $\nabla \Omega(\mu) \in \partial \Omega(\mu) \subset \mathbb{R}^k$. Additionally, if $\ell : \mathcal{M} \times \mathcal{Y} \rightarrow \bar{\mathbb{R}}$, then $\partial \Omega(\mu) \subset \mathbb{R}^k$, and if $\mu \in \text{relint } \mathcal{M}$, we have $\partial \Omega(\mu) \subset \mathbb{R}^k$. The loss is strictly proper if and only if the associated Ω is strictly convex.

With this theorem, we have an essentially complete analogy with Theorem 14.2.1. There are subtleties in the proof because the mapping from probabilities P to $\mathbb{E}_P[Y]$ can be many-to-one, necessitating some care in the calculations, and making infinite losses somewhat challenging. A few examples centered around ordinal regression illustrate the scenarios.

Example 14.2.16 (Ordinal regression, Example 14.2.14 continued): Let $Y \in \{0, 1, \dots, k\}$ be a value to be predicted, where the ordering on Y is important, as in ratings of items. In this case, the set $\mathcal{M} = \text{Conv}(\mathcal{Y}) = [0, k]$, and any strictly convex loss with domain $[0, k]$ gives rise to a proper loss via the construction $\ell_\Omega(\mu, y) = -\Omega(\mu) - \Omega'(\mu)(y - \mu)$. First, we take $\Omega(\mu) = \frac{1}{2}\mu^2$. This gives rise to a (modified) squared error

$$\ell_\Omega(\mu, y) = \frac{1}{2}(\mu - y)^2 - \frac{1}{2}y^2,$$

which is strictly convex and proper.

Other choices of Ω are possible. One natural choice is a variant of the negative binary entropy, and we define

$$\Omega(\mu) = (k - \mu) \log(k - \mu) + \mu \log \mu,$$

which is convex in $\mu \in [0, k]$, with $\Omega(\mu) = +\infty$ for $\mu > k$ or $\mu < 0$, while $\Omega(0) = \Omega(k) = k \log k$. We have $\Omega'(\mu) = \log \frac{\mu}{k - \mu}$, and so

$$\ell_\Omega(\mu, y) = -y \log \mu + (y - k) \log(k - \mu),$$

for $y \in \{0, \dots, k\}$. Here, however, note the importance of allowing infinite values in the loss ℓ when $\mu \rightarrow \{0, k\}$. \diamond

Proof One direction is, as in the previous cases, straightforward. Let ℓ have the given representation. Then for $\mu(P) = \mathbb{E}_P[Y]$,

$$\mathbb{E}_P[\ell(\mu, Y)] = -\Omega(\mu) - \langle \nabla \Omega(\mu), \mu(P) - \mu \rangle \geq -\Omega(\mu(P)) = \mathbb{E}_P[\ell(\mu(P), Y)],$$

and the inequality is strict if Ω is strictly convex.

The converse direction (from a proper loss to function Ω) is more subtle. We first give the argument in the case that the losses ℓ are finite-valued, so that $\ell(\mu, y) < \infty$ for each $\mu \in C$ and $y \in \mathcal{Y}$, deferring the proof of the general case to Section 14.5.1 as it yields little additional intuition. Let $\mathcal{Y} = \{y_1, \dots, y_m\} \subset \mathbb{R}^k$, and assume w.l.o.g. that the matrix $A = [y_1 \cdots y_m]$ with columns y_j has rank k (otherwise, we simply work in a subspace). We may identify \mathcal{P} with the probability simplex Δ_m , and then the mean mapping $\mu(p) = \sum_{i=1}^m p_i y_i$ for $p \in \mathbb{R}^m$ is surjective. Now for $\mu \in \mathbb{R}^k$ define

$$\Omega(\mu) := \inf_{p: \mu(p)=\mu} \sup_{\alpha} \{-\mathbb{E}_p[\ell(\alpha, Y)]\} \stackrel{(\star)}{=} \inf_{p \in \Delta_m} \{-\mathbb{E}_p[\ell(\mu, Y)] \mid \mu(p) = \mu\},$$

where the equality (\star) follows because ℓ is proper. The function Ω is closed convex, as it is the partial infimum of the closed convex function $p \mapsto -\mathbb{E}_p[\ell(\mu, Y)] + \mathbf{I}_{\Delta_m}(p)$, where we recall $\mathbf{I}_{\Delta_m}(p) = 0$ if $p \in \Delta_m$ and $+\infty$ otherwise (see Proposition B.3.11).

We compute $\partial\Omega(\mu)$ directly now. The infimum over p in the definition of $\Omega(\mu)$ is attained, as Δ_m is compact and $g(p) := -\mathbb{E}_p[\ell(\mu, Y)]$ is necessarily continuous in p satisfying $\mu(p) = \mu$, because regularity of the loss guarantees $\ell(\mu, y_i) \in \mathbb{R}$ whenever $p_i > 0$ is feasible in the mean mapping constraint $\mu(p) = \mu$. Moreover, it is immediate that

$$\nabla g(p) = \begin{bmatrix} -\ell(\mu, y_1) \\ \vdots \\ -\ell(\mu, y_m) \end{bmatrix} \in \mathbb{R}^m.$$

Let $p^*(\mu)$ be any p attaining the infimum. By Proposition B.3.29 on the subgradients of partial minimization, we thus obtain

$$\partial\Omega(\mu) = \left\{ s \in \mathbb{R}^k \mid y_i^T s = -\ell(\mu, y_i) \text{ for } i = 1, \dots, m \right\},$$

and moreover, this set is necessarily non-empty for all $\mu \in \text{relint } C = \{\mu(p) \mid p \succ 0, p \in \Delta_m\}$. Using this equality, we have

$$\begin{aligned} \ell(\mu, y) &= -\Omega(\mu) + \Omega(\mu) + \ell(\mu, y) = -\Omega(\mu) + \mathbb{E}_{p^*(\mu)}[-\ell(\mu, Y)] + \ell(\mu, y) \\ &= -\Omega(\mu) + \sum_{i=1}^m p_i^*(\mu) y_i^T s - y^T s \\ &= -\Omega(\mu) + \langle s, \mathbb{E}_{p^*(\mu)}[Y] - y \rangle = -\Omega(\mu) - \langle s, y - \mu \rangle \end{aligned}$$

for any $s \in \partial\Omega(\mu)$, as $\mathbb{E}_p[Y] = \mu(p) = \mu$ by construction.

Lastly, to obtain strict convexity of Ω , note that if $\mathbb{E}_p[Y] = \mu$, then we can use the representation

$$\mathbb{E}_p[\ell(\mu', Y)] - \mathbb{E}_p[\ell(\mu, Y)] = -\Omega(\mu') - \langle \nabla \Omega(\mu'), \mu - \mu' \rangle + \Omega(\mu) = D\Omega(\mu, \mu')$$

which is positive whenever $\mu \neq \mu'$ if and only if Ω is strictly convex. \square

14.3 From entropies to convex losses, arbitrary predictions, and link functions

Frequently, when we fit models, it is inconvenient to directly model or predict probabilities, that is, to minimize over probabilistic predictions. Instead, we often wish to fit some real-valued prediction and then transform it into a probabilistic prediction. This is perhaps most familiar from binary and multiclass logistic regression, where a *link function* transforms real-valued predictions into probabilistic predictions. For the binary logistic regression case with $Y \in \{-1, 1\}$, we assume that we predict a score $s \in \mathbb{R}$, where $s > 0$ indicates a prediction that Y is more likely to be 1 and $s < 0$ that it is more likely negative. The implied (modelled) probability that $Y = y$ is then

$$p(y | s) = \frac{1}{1 + \exp(-ys)} \quad \text{for } y \in \{-1, 1\}.$$

Similarly, for k -class classification problems, when using multiclass logistic regression, we predict a score vector $s \in \mathbb{R}^k$, where s_y indicates a score associated to one of the k potential class labels y ; this then implies the probabilities

$$p(y | s) = \frac{\exp(s_y)}{\sum_{i=1}^k \exp(s_i)} = \frac{1}{1 + \sum_{i \neq y} \exp(s_i - s_y)},$$

where we clearly have $\sum_y p(y | s) = 1$.

In binary and multiclass logistic regression, instead of directly minimizing negative log probabilities of error over the probability simplex (though one does this implicitly), instead we use surrogate logistic losses whose arguments can range over all of \mathbb{R} or \mathbb{R}^k . In the case of binary logistic regression with $y \in \{-1, 1\}$, this is

$$\varphi(s, y) = \log(1 + \exp(-sy)),$$

while in the multiclass case we use the multiclass logistic loss

$$\varphi(s, y) = -s_y + \log \left(\sum_{i=1}^k \exp(s_i) \right) = \log \left(1 + \sum_{i \neq y} \exp(s_i - s_y) \right).$$

Note that for each of these, we have a direct relationship between the probabilistic predictions and derivatives of φ . In the binary logistic regression case, we have

$$p(y | s) = 1 + \frac{\partial}{\partial s} \varphi(s, y) = 1 - \frac{1}{1 + \exp(ys)} = \frac{1}{1 + \exp(-ys)},$$

while in the multiclass case we similarly have

$$p(y | s) = 1 + \frac{\partial}{\partial s_y} \varphi(s, y) = \frac{\exp(s_y)}{\sum_{i=1}^k \exp(s_i)}.$$

14.3.1 Convex conjugate linkages

These dualities turn out to hold in substantially more generality, and they are the key to transforming proper losses (as applied on probabilities) into proper surrogate losses that apply directly to real-valued scores and which are convex in their arguments, allowing us to bring the tools of

convex optimization to bear on actually fitting predictive models. We work in the general setting of Section 14.2.3 of losses for vector-valued y where $\mathcal{Y} \subset \mathbb{R}^k$, so that instead of predicting probability distributions on Y itself we predict elements μ of the set $\{\mathbb{E}_P[Y]\} = \text{Conv}(\mathcal{Y})$, and let ℓ be a strictly proper loss. Theorems 14.2.1 and 14.2.15 demonstrate that if the loss ℓ is proper, there exists a (negative) generalized entropy, which in the case of Theorem 14.2.1 is $\Omega(p) = \sup_q \{-\mathbb{E}_p[\ell(q, Y)]\}$, for which

$$\ell(\mu, y) = -\Omega(\mu) - \langle \nabla \Omega(\mu), y - \mu \rangle.$$

Note that Ω is always a *closed* convex function, meaning that it is lower semicontinuous or that its epigraph $\text{epi } \Omega = \{(\mu, t) \mid \Omega(\mu) \leq t\}$ is closed.

Let us suppose temporarily that we have *any* such entropy. Recalling the convex conjugate (14.1.3), the negative generalized entropy Ω is closed convex, and so its conjugate $\Omega^*(s) = \sup\{\langle s, \mu \rangle - \Omega(\mu)\}$ satisfies $\Omega^{**}(\mu) = \Omega(\mu)$. In particular, if we define the *surrogate loss*

$$\varphi(s, y) := \Omega^*(s) - \langle s, y \rangle,$$

which is defined for all $s \in \mathbb{R}^k$ (instead of $\text{Conv}(\mathcal{Y})$), then

$$\mathbb{E}_P[\varphi(s, Y)] = \Omega^*(s) - \langle s, \mathbb{E}_P[Y] \rangle = \Omega^*(s) - \langle s, \mu(P) \rangle$$

for the mean mapping $\mu(P) = \mathbb{E}_P[Y]$. Moreover,

$$\inf_s \mathbb{E}_P[\varphi(s, Y)] = \inf_s \{\Omega^*(s) - \langle s, \mu(P) \rangle\} = -\Omega^{**}(\mu(P)) = -\Omega(\mu(P)),$$

and so it generates the same negative entropy as the original loss ℓ , as

$$\inf_\mu \mathbb{E}_P[\ell(\mu, Y)] = \inf_\mu \{-\Omega(\mu) - \langle \nabla \Omega(\mu), \mu(P) - \mu \rangle\} = -\Omega(\mu(P)).$$

This identification of (generalized) entropies will underpin much of our development of the consistency of losses in sections to come. For now, we content ourselves with addressing how to understand propriety of the surrogate loss φ and how to transform predictions $s \in \mathbb{R}^k$ into probabilistic predictions μ .

The key will be to consider what we term *convex-conjugate-linkages*, or *conjugate linkages* for short. Recall the duality relationships (14.1.4) from the Fenchel-Young inequality we present in the convexity primer in Section 14.1.1. The negative generalized entropy Ω is convex, and the dualities associated with its conjugate $\Omega^*(s) = \sup_\mu \{\langle s, \mu \rangle - \Omega(\mu)\}$ will form the basis of our transformations. We first give a somewhat heuristic presentation, as the intuition is important (but details to make things precise can be a bit tedious). Essentially, we require that Ω^* and Ω are continuously differentiable, in which case we have

$$\nabla \Omega(\mu) = s \text{ if and only if } \nabla \Omega^*(s) = \mu \text{ if and only if } \Omega^*(s) + \Omega(\mu) = \langle s, \mu \rangle$$

by the Fenchel-Young inequalities (14.1.4). That is, the gradient $\nabla \Omega^*$ of the conjugate transforms a score vector $s \in \mathbb{R}^k$ into elements c to predict \mathcal{Y} : we transform s into a prediction μ via the *conjugate link* function

$$\text{pred}_\Omega(s) = \arg\max_\mu \{\langle s, \mu \rangle - \Omega(\mu)\} = \nabla \Omega^*(s) = (\nabla \Omega)^{-1}(s), \quad (14.3.1)$$

which finds the μ that best trades having maximal “entropy” $-\Omega(\mu)$, or uncertainty, with alignment with the scores $\langle s, \mu \rangle$.

With this, it is then natural to consider the function substituting the prediction $\mu = \text{pred}_\Omega(s)$ into $\ell(\mu, y)$, and so we consider

$$\ell(\text{pred}_\Omega(s), y).$$

Immediately, if $\mu = \text{pred}_\Omega(s) = \nabla\Omega^*(s)$, we have $s = \nabla\Omega(\mu)$ by construction (or the Fenchel-Young inequality (14.1.4)), and so $\Omega(\mu) = \langle s, \mu \rangle - \Omega^*(s)$ for this particular pair (s, μ) , and $\nabla\Omega(\mu) = \nabla\Omega(\nabla\Omega^*(s)) = s$ because $\nabla\Omega$ and $\nabla\Omega^*$ are inverses. Substituting, we obtain

$$\begin{aligned} \ell(\text{pred}_\Omega(s), y) &= -\Omega(\text{pred}_\Omega(s)) - \langle \nabla\Omega(\text{pred}_\Omega(s)), y - \text{pred}_\Omega(s) \rangle = -\Omega(\mu) - \langle s, y - \mu \rangle \\ &= \Omega^*(s) - \langle s, \mu \rangle - \langle s, y - \mu \rangle, \end{aligned}$$

that is, we have recovered the surrogate

$$\varphi(s, y) = \Omega^*(s) - \langle s, y \rangle. \quad (14.3.2)$$

The surrogate loss (14.3.2) constructed from the negative entropy Ω is the key transformation of the loss ℓ into a convex loss, and (no matter the properties of ℓ) is always convex.

As we have already demonstrated, the construction (14.3.2) is more general than we have presented; certainly, Ω^* is always convex, and so φ is always convex in s . Moreover, if Y has expectation $\mathbb{E}[Y] = \mu$, then

$$\inf_s \mathbb{E}[\varphi(s, \mu)] = \inf_s \{\Omega^*(s) - \langle s, \mu \rangle\} = -\Omega(\mu)$$

by conjugate duality, so the surrogate φ always recovers the negative entropy Ω ; without some type of differentiability conditions, however, the construction of the prediction mapping pred_Ω requires more care. Sections 16.3, 16.4, and 16.5 more deeply investigate these connections.

All that remains is to give more precise conditions under which the prediction (14.3.1) is always unique and exists for all possible score vectors $s \in \mathbb{R}^k$. To that end, we make the following definition.

Definition 14.4. *Let $\Omega : \mathbb{R}^k \rightarrow \overline{\mathbb{R}}$. Then Ω is a Legendre negative entropy if it is strictly convex, continuously differentiable on the interior of its domain, and*

$$\|\nabla\Omega(\mu)\| \rightarrow \infty \quad \text{if either} \quad \begin{cases} \mu \rightarrow \text{bd dom } \Omega & \text{or} \\ \|\mu\| \rightarrow \infty. \end{cases} \quad (14.3.3)$$

This is precisely the condition we require to make each step in the development of the surrogate (14.3.2) airtight; as a corollary to Theorem C.2.9 in the appendices, we have the following.

Corollary 14.3.1. *Let Ω be a Legendre negative entropy. Then the conjugate link prediction (14.3.1) is unique and exists for all $s \in \mathbb{R}^k$. In particular, the conjugate Ω^* is strictly convex, continuously differentiable, satisfies $\text{dom } \Omega^* = \mathbb{R}^k$, and $\nabla\Omega^* = (\nabla\Omega)^{-1}$.*

With this corollary in place, we can then give a theorem showing the equivalence of the strictly proper loss ℓ and its surrogate.

Theorem 14.3.2. *Let $\ell : \mathcal{M} \times \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ be the strictly proper loss associated with the Legendre negative entropy Ω . Then*

$$\ell(\text{pred}_\Omega(s), y) = \varphi(s, y) := \Omega^*(s) - \langle s, y \rangle.$$

Moreover, the convex surrogate φ satisfies the consistency that if

$$\mathbb{E}[\varphi(s_n, Y)] \rightarrow \inf_s \mathbb{E}[\varphi(s, Y)]$$

then $\mu_n = \text{pred}_\Omega(s_n)$ satisfies

$$\mathbb{E}[\ell(\mu_n, Y)] \rightarrow \inf_\mu \mathbb{E}[\ell(\mu, Y)].$$

Proof The first equality we have already demonstrated. For the minimization claim, we note that if $\mu = \mathbb{E}[Y]$, then $\mathbb{E}[\varphi(s, Y)] = \Omega^*(s) - \langle \mu, s \rangle$ and $\inf_s \{\Omega^*(s) - \langle \mu, s \rangle\} = -\Omega(\mu)$. Strict propriety of ℓ then gives $\inf_{\mu'} \mathbb{E}[\ell(\mu', Y)] = -\Omega(\mu)$. \square

Said differently, the surrogate φ is consistent with the loss ℓ and (strictly) proper, in that if s minimizes $\mathbb{E}[\varphi(s, Y)]$, then $\text{pred}_\Omega(s)$ minimizes $\mathbb{E}[\ell(\mu, Y)]$. The statement in terms of limits is necessary, however, as simple examples show, because with some link functions it is in fact impossible to achieve the extreme points of $\text{Conv}(\mathcal{Y})$, as in logistic regression. We provide a few example applications (and non-applications) of Theorem 14.3.2. For the first, let us consider binary logistic regression.

Example 14.3.3 (Binary logistic regression): For a label $Y \in \{0, 1\}$ and predictions $p \in [0, 1]$, take the generalized negative entropy

$$\Omega(p) = p \log p + (1 - p) \log(1 - p).$$

By inspection, $\text{dom } \Omega = [0, 1]$, and $\Omega'(p) = \log \frac{p}{1-p}$ satisfies $|\Omega'(p)| \rightarrow \infty$ as $p \rightarrow \{0, 1\}$. For $s \in \mathbb{R}$, the conjugate is

$$\Omega^*(s) = \sup_p \{sp - p \log p - (1 - p) \log(1 - p)\} = \log(1 + e^s),$$

where the supremum is achieved by $p = \text{pred}_\Omega(s) = \frac{e^s}{1+e^s}$. Then we have

$$\varphi(s, y) = \log(1 + e^s) - sy = -\log p(y | s),$$

where $p(y | s) = \frac{e^{ys}}{1+e^s}$ is the binary logistic probability of the label $y \in \{0, 1\}$.

For the induced loss $\ell(p, y) = -y \log p - (1 - y) \log(1 - p)$ (the log loss), if $\mathbb{P}(Y = 1) = 1$, then $p = 1$ minimizes $\mathbb{E}[\ell(p, Y)]$. Similarly, if $\mathbb{P}(Y = 0) = 1$, then $p = 0$ minimizes $\mathbb{E}[\ell(p, Y)]$. Neither of these is achievable by a finite $\hat{\ell}$ in $p(y | s) = \frac{e^{ys}}{1+e^s}$, showing how the limiting argument in Theorem 14.3.2 is necessary. \diamond

The next example shows that we sometimes need to elaborate the setting of Theorem 14.3.2 to deal with constraints.

Example 14.3.4 (Multiclass logistic regression): Identify the set $\mathcal{Y} = \{e_1, \dots, e_k\}$ with the k standard basis vectors, and for $p \in \Delta_k = \{p \in \mathbb{R}_+^k \mid \mathbf{1}^T p = 1\}$, consider the negative entropy

$$\Omega(p) = \sum_{y=1}^k p_y \log p_y.$$

This function is strictly convex and of Legendre type for the positive orthant \mathbb{R}_+^k but *not* for Δ_k . Shortly, we shall allow linear constraints on the predictions to address this shortcoming.

As an alternative, take $\mathcal{Y} = \{0, e_1, \dots, e_{k-1}\}$, so that $\text{Conv}(\mathcal{Y}) = \{p \in \mathbb{R}_+^{k-1} \mid \mathbf{1}^T p \leq 1\}$, which has an interior and so more easily admits a conjugate duality relationship. In this case, the negative entropy-type function

$$\Omega(p) = \sum_{y=1}^{k-1} p_y \log p_y + (1 - \mathbf{1}^T p) \log(1 - \mathbf{1}^T p) \quad (14.3.4)$$

is of Legendre type. A calculation for $s \in \mathbb{R}^{k-1}$ yields

$$\Omega^*(s) = \log \left(1 + \sum_{y=1}^{k-1} e^{s_y} \right),$$

with

$$\text{pred}_\Omega(s) = \left(\frac{e^{s_1}}{1 + \sum_{j=1}^{k-1} e^{s_j}}, \dots, \frac{e^{s_{k-1}}}{1 + \sum_{j=1}^{k-1} e^{s_j}} \right).$$

Letting p denote the entries of this vector, we can then assign a probability to class k via $p_k = 1 - \sum_{j=1}^{k-1} p_j$. \diamond

In Section 14.4 we revisit exponential families in the (proper) loss minimization framework we have thus far developed, which gives some additional perspective on these problems.

14.3.2 Convex conjugate linkages with affine constraints

As Example 14.3.4 shows, in some cases a “natural” formulation fails to satisfy the desiderata of our link functions. Accordingly, we make a slight modification to the Legendre type (14.3.3) negative entropy h to allow for *affine* constraints, which still allows us to develop the precise convexity dualities with proper losses we require. Continuing to work in the scenario in which $\mathcal{Y} \subset \mathbb{R}^k$, suppose now that the affine hull

$$\mathcal{A} = \text{aff}(\mathcal{Y}) := \left\{ \sum_{j=1}^m \alpha_j y_j \mid y_j \in \mathcal{Y}, \alpha^T \mathbf{1} = 1, m \in \mathbb{N} \right\}$$

is a proper subspace of \mathbb{R}^k . The key motivating example here is the “failure” case of Example 14.3.4 on multiclass logistic regression, where $\mathcal{Y} = \{e_1, \dots, e_k\}$, whose affine hull is exactly those vectors $p \in \mathbb{R}^k$ satisfying $\langle p, \mathbf{1} \rangle = 1$. Naturally, in this case we wish to predict probabilities, and so given a score vector $s \in \mathbb{R}^k$ and using the negative entropy $\Omega(p) = \sum_{y=1}^k p_y \log p_y$, we let

$$\text{pred}(s) = \underset{p}{\text{argmin}} \{ \Omega(p) - \langle s, p \rangle \mid \mathbf{1}^T p = 1 \} = \left[\frac{e^{s_y}}{\sum_{j=1}^k e^{s_j}} \right]_{y=1}^k.$$

Generalizing this approach to arbitrary regularizers Ω , we modify the prediction (14.3.1) to be

$$\text{pred}_{\Omega, \mathcal{A}}(s) = \underset{\mu \in \mathcal{A}}{\text{argmax}} \{ \langle s, \mu \rangle - \Omega(\mu) \}.$$

Then for the loss $\ell(\mu, y) = -\Omega(\mu) - \langle \nabla \Omega(\mu), y - \mu \rangle$ associated with the negative entropy Ω , we define the surrogate

$$\varphi(s, y) := \ell(\text{pred}_{\Omega, \mathcal{A}}(s), y).$$

Perhaps remarkably, this construction still yields a well-defined convex loss with the same consistency properties as those in Theorem 14.3.2. Indeed, defining

$$\Omega_{\mathcal{A}}(\mu) = \Omega(\mu) + \mathbf{I}_{\mathcal{A}}(\mu)$$

and the associated conjugate $\Omega_{\mathcal{A}}^*(s) = \sup\{\langle s, \mu \rangle - \Omega(\mu) \mid \mu \in \mathcal{A}\}$, we have the following theorem.

Theorem 14.3.5. *Let $\ell : \mathcal{M} \times \mathcal{Y} \rightarrow \bar{\mathbb{R}}$ be the strictly proper loss associated with the Legendre negative entropy Ω and $\mathcal{A} = \text{aff}(\mathcal{Y})$ be the affine hull of \mathcal{Y} . Then*

$$\varphi(s, y) := \ell(\text{pred}_{\Omega, \mathcal{A}}(s), y) = \Omega_{\mathcal{A}}^*(s) - \langle s, y \rangle.$$

Moreover, the convex surrogate φ satisfies the consistency that if

$$\mathbb{E}[\varphi(s_n, Y)] \rightarrow \inf_s \mathbb{E}[\varphi(s, Y)]$$

then $\mu_n = \text{pred}_{\Omega, \mathcal{A}}(s_n)$ satisfies

$$\mathbb{E}[\ell(\mu_n, Y)] \rightarrow \inf_{\mu} \mathbb{E}[\ell(\mu, Y)].$$

We return to proving the theorem presently, focusing here on how it applies to Example 14.3.4.

Example 14.3.6 (Multiclass logistic regression): Consider Example 14.3.4, where we identify $\mathcal{Y} = \{e_1, \dots, e_k\} \subset \mathbb{R}^k$, which has affine hull $\mathcal{A} = \{p \in \mathbb{R}^k \mid \langle \mathbf{1}, p \rangle = 1\}$. Then taking $\Omega(p) = \sum_{k=1}^k p_k \log p_k$, a calculation with a Lagrangian shows that

$$\text{pred}_{\Omega, \mathcal{A}}(s) = \underset{p \in \Delta_k}{\text{argmin}} \{-\langle s, p \rangle + \Omega(p)\} = \left[\frac{e^{s_y}}{\sum_{j=1}^k e^{s_j}} \right].$$

In turn, this gives surrogate logistic loss

$$\varphi(s, y) = \log \left(\sum_{j=1}^k e^{s_j - s_y} \right).$$

Notably, the logistic loss is *not* strictly convex, as $\varphi(s + t\mathbf{1}, y) = \varphi(s, y)$ for $t \in \mathbb{R}$. If Y is a multinomial random variable with $\mathbb{P}(Y = e_y) = p_y$, then by another calculation, the vector with entries

$$s_y^* = \log p_y$$

minimizes $\mathbb{E}[\varphi(s, Y)]$, which in turn gives $\text{pred}_{\Omega, \mathcal{A}}(s^*) = p$, maintaining propriety. \diamond

Proof of Theorem 14.3.5

Before proving the theorem proper, we show how the key identity that $s = \nabla\Omega(\mu)$ we use to develop equality (14.3.2) generalizes in the presence of the affine constraint. The function $h_{\mathcal{A}}$ is strictly convex on its domain $\text{dom } h \cap \mathcal{A}$, and moreover, $\nabla h_{\mathcal{A}}^*$ exists and is continuous. The following corollary (a consequence of Corollary C.2.12 in Appendix C.2) extends Corollary 14.3.1 and allows us to address equality (14.3.2).

Corollary 14.3.7. *The conjugate $\Omega_{\mathcal{A}}^*$ is continuously differentiable with $\text{dom } \Omega_{\mathcal{A}}^* = \mathbb{R}^k$, and if $\mu = \nabla\Omega_{\mathcal{A}}^*(s)$, then $\mu \in \text{int dom } \Omega$ and*

$$\nabla\Omega(\mu) = s + v$$

for some vector v normal to \mathcal{A} , that is, a vector $v \in \mathbb{R}^k$ satisfying $\langle v, \mu_0 - \mu_1 \rangle = 0$ for all $\mu_0, \mu_1 \in \mathcal{A}$.

While the proof of the corollary requires some care to make precise, a sketch can give intuition.

Sketch of Proof Because Ω is strictly convex and its derivatives $\nabla\Omega(\mu)$ explode as $\mu \rightarrow \text{bd dom } \Omega$, the minimizer of $-\langle s, \mu \rangle + \Omega(\mu)$ over $\mu \in \mathcal{A}$ exists and is unique. Let $\mathcal{A} = \{\mu \mid A\mu = b\}$ for shorthand, where $A \in \mathbb{R}^{n \times k}$ for some $n < k$. Then introducing Lagrange multiplier $w \in \mathbb{R}^n$ for the constraint $\mu \in \mathcal{A}$, the Lagrangian for finding $\text{pred}_{\Omega, \mathcal{A}}(s) = \text{argmin}_{\mu} \{\Omega(\mu) - \langle s, \mu \rangle \mid \mu \in \mathcal{A}\}$ is

$$\mathcal{L}(\mu, w) = \Omega(\mu) - \langle s, \mu \rangle + w^T(A\mu - b).$$

Minimizing out μ by setting $\nabla_{\mu}\mathcal{L}(\mu, w) = 0$, we obtain

$$\nabla\Omega(\mu) - s + A^T w = 0.$$

But if $\mu_0, \mu_1 \in \mathcal{A}$, then $v = A^T w$ satisfies $\langle v, \mu_0 - \mu_1 \rangle = w^T A(\mu_0 - \mu_1) = w^T(b - b) = 0$, so that v is normal to \mathcal{A} . \square

Finally, we return to prove the theorem. Take any vector $s \in \mathbb{R}^k$. Then because $\text{pred}_{\Omega, \mathcal{A}}(s) = \nabla\Omega_{\mathcal{A}}^*(s)$, we have

$$\varphi(s, y) = \ell(\text{pred}_{\Omega, \mathcal{A}}(s), y) = -\Omega(\nabla\Omega_{\mathcal{A}}^*(s)) - \langle \nabla\Omega(\nabla\Omega_{\mathcal{A}}^*(s)), y - \nabla\Omega_{\mathcal{A}}^*(s) \rangle.$$

As $\nabla\Omega_{\mathcal{A}}^*(s) \in \mathcal{A}$ and using the shorthand $\mu = \nabla\Omega_{\mathcal{A}}^*(s) \in \mathcal{A}$, we have $\nabla\Omega(\mu) = s + v$ for some v normal to \mathcal{A} . Moreover, $\Omega(\mu) = \Omega_{\mathcal{A}}(\mu)$, and so the Fenchel-Young inequality (14.1.4) guarantees $-\Omega_{\mathcal{A}}(\mu) = \Omega_{\mathcal{A}}^*(s) - \langle s, \mu \rangle$. Substituting in the expression for φ , we obtain

$$\begin{aligned} \varphi(s, y) &= \Omega_{\mathcal{A}}^*(s) - \langle s, \mu \rangle - \langle s + v, y - \mu \rangle \\ &= \Omega_{\mathcal{A}}^*(s) - \langle s, \mu \rangle + \langle s, \mu - y \rangle = \Omega_{\mathcal{A}}^*(s) - \langle s, y \rangle \end{aligned}$$

where the second equality follows because $v \perp \mu - y$.

For the consistency argument, let $\mu_n = \text{pred}_{\Omega, \mathcal{A}}(s_n)$. Then $\mathbb{E}[\ell(\mu_n, Y)] = \mathbb{E}[\varphi(s_n, Y)]$ and if $\mu = \mathbb{E}[Y]$, then $\mathbb{E}[\varphi(s, Y)] = \Omega_{\mathcal{A}}^*(s) - \langle \mu, s \rangle$ and $\inf_s \mathbb{E}[\varphi(s, Y)] = -\Omega_{\mathcal{A}}(\mu) = -\Omega(\mu)$. Strict propriety of ℓ gives $\inf_{\mu'} \mathbb{E}[\ell(\mu', Y)] = -\Omega(\mu)$.

14.4 Exponential families, maximum entropy, and log loss

Realistically, making predictions using an arbitrary distribution P on an arbitrary space \mathcal{X} is statistically infeasible: we could never collect enough data to accurately model complex phenomena without any assumptions on P . Accordingly, we may seek more tractable models to make predictions feasible, and we can then investigate the consequences of moving from the entire family \mathcal{P} of distributions on \mathcal{X} to smaller families is. A particularly important class of distributions, which allows us to study these questions in great detail, are the exponential families from Chapter 3; here, we investigate them in the framework that we have developed for proper losses.

Let $\{P_\theta\}$ be a regular exponential family indexed by θ on a space \mathcal{X} with sufficient statistic $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$, where for a base measure ν on \mathcal{X} , P_θ has density

$$p_\theta(x) = \exp(\langle \theta, \phi(x) \rangle - A(\theta))$$

with respect to ν , where $A(\theta) = \log \int e^{\langle \theta, \phi(x) \rangle} d\nu(x)$ is the log partition function. (Recall that regularity means that the domain

$$\Theta := \text{dom } A = \{\theta \mid A(\theta) < \infty\}$$

is open, as in Definition 3.1). Consider the log loss $-\log p_\theta(x)$, which we suggestively denote with the surrogate φ as a function of θ ,

$$\varphi(\theta, x) := -\log p_\theta(x) = A(\theta) - \langle \theta, \phi(x) \rangle.$$

Proposition 3.2.1 guarantees this is always convex in θ because the log partition function is convex, and it is \mathcal{C}^∞ (Proposition 3.2.2). While the log loss $-\log p(x)$ is proper, the exponential family $\{P_\theta\}$ can capture only a subset of the distributions on \mathcal{X} .

The mean mapping $\mu(P) := \mathbb{E}_P[\phi(X)] \in \mathbb{R}^d$ will be of central importance to the development of proper losses, exponential families, and the duality relationships between maximum likelihood and entropy that we explore here. Accordingly, throughout this section we let

$$\mathcal{P} := \{\text{distributions } P \ll \nu\} = \{\text{distributions } P \text{ with a density } p \text{ w.r.t. } \nu\}$$

be the collection of distributions with densities with respect to ν (as P_θ by definition has), and we define the set of potential *mean parameters*

$$\mathcal{M} := \{\mu(P) = \mathbb{E}_P[\phi(X)] \in \mathbb{R}^d \mid P \ll \nu\} = \{\mu(P) \mid P \in \mathcal{P}\}. \quad (14.4.1)$$

Now, for any distribution $P \in \mathcal{P}$ with mean vector $\mu = \mu(P)$, the associated generalized negative entropy is

$$\Omega(\mu) := \sup_{\theta} \{-\mathbb{E}_P[\varphi(\theta, X)]\} = \sup_{\theta} \{\langle \theta, \mu(P) \rangle - A(\theta)\} = A^*(\mu),$$

the convex conjugate of A . At this point, the centrality of the duality relationships (via gradients ∇A and ∇A^*) between Θ and \mathcal{M} to fitting and modeling should come as no surprise, and so we elucidate a few of the main properties. Because $\nabla A(\theta) = \mathbb{E}_\theta[\phi(X)]$ in the exponential family, we immediately see that

$$\nabla A(\Theta) := \{\nabla A(\theta)\}_{\theta \in \Theta} \subset \mathcal{M}.$$

Recalling the duality relationship (14.1.4) that

$$\theta \in \partial A^*(\mu) \text{ if and only if } \nabla A(\theta) = \mu,$$

we can say much more.

Proposition 14.4.1. *Let $\mathcal{M}^\circ = \text{relint } \mathcal{M}$. Then $\nabla A(\Theta) = \mathcal{M}^\circ$. Additionally:*

- (i) *If the family is minimal, then \mathcal{M} has non-empty interior and Ω is continuously differentiable on \mathcal{M}° , with $\theta = \nabla \Omega(\mu)$ if and only if $\nabla A(\theta) = \mu$.*
- (ii) *If the family is non-minimal, then Ω is continuously differentiable relative to $\text{aff}(\mathcal{M})$, meaning that there exists a continuous mapping $\nabla \Omega(\mu) \in \Theta$ such that for all $\mu \in \mathcal{M}^\circ$,*

$$\partial \Omega(\mu) = \left\{ \nabla \Omega(\mu) + \text{aff}(\mathcal{M})^\perp \right\}.$$

Moreover, $\Theta = \Theta + \text{aff}(\mathcal{M})^\perp$.

The proof of the proposition relies on the more sophisticated duality theory we develop in Appendices B and C, so we defer it to Section 14.5.2.

We can summarize the proposition by considering minimizers and maximizers: suppose we wish to choose θ to minimize

$$\mathbb{E}_P[\varphi(\theta, X)] = \mathbb{E}_P[-\log p_\theta(X)] = A(\theta) - \langle \mu(P), \theta \rangle.$$

Then so long as the distribution P is not extremal in that $\mu(P) = \mathbb{E}_P[\phi(X)] \in \text{relint } \mathcal{M}$, there exists a parameter $\theta(P)$, unique up to translation in the subspace perpendicular to $\text{aff}(\mathcal{M})$, for which

$$\theta(P) \in \underset{\theta}{\text{argmin}} \mathbb{E}_P[\varphi(\theta, X)] = \underset{\theta}{\text{argmin}} \{A(\theta) - \langle \mu(P), \theta \rangle\}.$$

Moreover, this parameter satisfies the mean matching condition

$$\nabla A(\theta(P)) = \mu(P),$$

which is of course sufficient to be a minimizer of the expected log loss. As the statements in the proposition evidence, calculations become more challenging when we must perform them all in an affine subspace, though sometimes this care is unavoidable.

Example 14.4.2 (Gaussian estimation): Assume we fit a distribution assuming X has a Gaussian distribution with mean μ and covariance $\Sigma \succ 0$, both to be estimated. Performing the transformation to the exponential family form with precision $K = \Sigma^{-1}$ and $\theta = \Sigma^{-1}\mu$, we have

$$p_{\theta, K}(x) = \exp \left(\langle \theta, x \rangle - \frac{1}{2} \langle xx^\top, K \rangle - A(\theta, K) \right) \quad \text{for} \quad A(\theta, K) = \frac{1}{2} \theta^\top K^{-1} \theta - \frac{1}{2} \log \det(2\pi K).$$

The log partition function has gradients

$$\nabla_\theta A(\theta, K) = K^{-1} \theta \quad \text{and} \quad \nabla_K A(\theta, K) = -\frac{1}{2} K^{-1} \theta \theta^\top K^{-1} - \frac{1}{2} K^{-1}.$$

Matching moments for a distribution P with second moment matrix $M = \mathbb{E}[XX^\top] \succ 0$ and mean $\mathbb{E}[X]$, we obtain

$$\mathbb{E}[X] = K^{-1} \theta \quad \text{and} \quad M = K^{-1} \theta \theta^\top K^{-1} + K^{-1}.$$

Setting $\theta = K \mathbb{E}[X]$ and noting that $M = \text{Cov}(X) - \mathbb{E}[X] \mathbb{E}[X]^\top$, we solve $M = \mathbb{E}[X] \mathbb{E}[X]^\top + K^{-1}$ by setting $K^{-1} = \text{Cov}(X)$.

When $\text{Cov}(X) \not\asymp 0$, the solution $K = \text{Cov}(X)^{-1}$ does not exist, so we must rely instead on part (ii) of Proposition 14.4.1. With some care, one may check that we can work in the subspace spanned by the eigenvectors of $\text{Cov}(X)$, that is, if $\text{Cov}(X) = U\Lambda U^\top$ and $U \in \mathbb{R}^{d \times k}$, the collection of symmetric matrices K whose column space belongs to $\text{span}(U)$. Then the pseudo-inverse $K = \text{Cov}(X)^\dagger$ is the appropriate solution, and it recovers the covariance $\Sigma = K^\dagger = \text{Cov}(X) \succeq 0$. \diamond

Finally, let us give a last result that shows the duality relationships between the negative generalized entropy $\Omega(\mu)$ and log partition A , which allows us to also capture a few of the nuances of minimization of the surrogate log loss $\varphi(\theta, x) = -\log p_\theta(x)$ when we encounter distributions P for which the mean mapping $\mu(P)$ is on the boundary of \mathcal{M} or even outside it.

Proposition 14.4.3. *Let $\{P_\theta\}$ be a regular exponential family with log partition $A(\theta)$ with domain Θ , and let \mathcal{M} be the associated mean parameter space with relative interior $\mathcal{M}^\circ = \text{relint } \mathcal{M}$. Let $\Omega(\mu) = A^*(\mu)$ be the associated negative generalized entropy. Then*

- (i) $A(\theta) = \Omega^*(\theta) = A^{**}(\theta)$ for all θ .
- (ii) If $\mu \in \mathcal{M}^\circ$, there exists $\theta(\mu) \in \Theta$ such that the negative entropy satisfies $\Omega(\mu) = A^*(\mu) = \langle \theta(\mu), \mu \rangle - A(\theta(\mu)) < \infty$. If $\mu \notin \text{cl } \mathcal{M}$, then $\Omega(\mu) = +\infty$.
- (iii) If $\mu \in \text{bd } \mathcal{M} = \text{cl } \mathcal{M} \setminus \mathcal{M}^\circ$, then for any $\mu_0 \in \mathcal{M}^\circ$, $\Omega(\mu) = \lim_{t \rightarrow 0} \Omega(t\mu_0 + (1-t)\mu)$, and there exist $\theta_t \in \Theta$ with

$$\nabla A(\theta_t) = t\mu_0 + (1-t)\mu \quad \text{and} \quad \lim_{t \rightarrow 0} \{A(\theta_t) - \langle \mu, \theta_t \rangle\} = \inf_{\theta} \{A(\theta) - \langle \mu, \theta \rangle\}.$$

In particular, there exist sequences of dual pairs (μ_n, θ_n) with $\mu_n \in \mathcal{M}^\circ$ and $\theta_n \in \Theta$ satisfying $\mu_n = \nabla A(\theta_n)$, $\mu_n \rightarrow \mu$, $\Omega(\mu_n) \rightarrow \Omega(\mu)$, and $A(\theta_n) - \langle \mu, \theta_n \rangle \rightarrow \inf_{\theta} \{A(\theta) - \langle \mu, \theta \rangle\}$.

See Section 14.5.2 for the deferred proof.

While the statement of Proposition 14.4.3 is complex, considering minimizers of $\mathbb{E}[\varphi(\theta, X)]$ can give some understanding. If P is a distribution such that $\mu(P) \in \mathcal{M}^\circ$, then there exists a parameter $\theta(P)$ minimizing $\mathbb{E}_P[\varphi(\theta, X)]$. If $\mu(P) \in \text{bd } \mathcal{M}$, then either there exists a minimizer $\theta(P)$ of the loss, or there is a sequence of points θ_n such that

$$\mathbb{E}_P[\varphi(\theta_n, X)] \rightarrow \inf_{\theta} \mathbb{E}_P[\varphi(\theta, X)] = -\Omega(\mu(P)), \quad \text{and} \quad \mu(P_{\theta_n}) \rightarrow \mu(P),$$

so that they asymptotically satisfy the mean identity. Finally, if $\mu(P) \notin \text{cl } \mathcal{M}$, then $\inf_{\theta} \mathbb{E}[\varphi(\theta, X)] = -\infty$, making the choice of exponential family model poor, as it cannot capture the mean parameters.

14.4.1 Maximizing entropy

As we have seen, our notion of generalized entropies as the minimal values of expected losses can recapture the classical entropy $H(P) = -\sum_x p(x) \log p(x)$ when P has a probability mass function p , as in the case of multiclass prediction. For exponential family models, this connection goes much further, and the negative generalized entropy $\Omega(\mu)$ for $\mu \in \mathcal{M}$ coincides with a more general notion of entropy known as the *Shannon entropy*. We begin with the definition:

Definition 14.5. *Let ν be a base measure on \mathcal{X} and assume P has density p with respect to ν . Then the Shannon entropy of P is*

$$H(P) = - \int p(x) \log p(x) d\nu(x).$$

For a distribution P with probability mass function p , the base measure ν is counting measure, yielding the classical entropy $H(P) = -\sum_x p(x) \log p(x)$, while for a distribution P with density p (for Lebesgue measure ν , so that $d\nu(x) = dx$ for $x \in \mathbb{R}^d$), we recover the differential entropy $H(P) = -\int p(x) \log p(x) dx$.

Example 14.4.4: Let P be the uniform distribution on $[0, a]$. Then the differential entropy $H(P) = -\log(1/a) = \log a$. \diamond

Example 14.4.5: Let P be the normal distribution $\mathcal{N}(\mu, \Sigma)$ and ν be Lebesgue measure. Then

$$H(P) = \frac{1}{2} \log(\det(2\pi\Sigma)) + \frac{1}{2} \mathbb{E}[(X - \mu)^\top \Sigma^{-1} (X - \mu)] = \frac{d}{2} \log(2\pi e) + \frac{1}{2} \log \det(\Sigma).$$

because $p(x) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1} (x - \mu))$. \diamond

For exponential families, the log partition determines the Shannon entropy directly, highlighting that $-h$ is indeed a familiar entropy-like object.

Proposition 14.4.6. Let $\{P_\theta\}$ be a regular exponential family with respect to the base measure ν . Then for any $\theta \in \Theta$,

$$H(P_\theta) = -h(\mu(P_\theta)) = A(\theta) - \langle \mu(P_\theta), \theta \rangle,$$

where $\Omega(\mu) = \sup\{\langle \mu, \theta \rangle - A(\theta)\} = A^*(\mu)$.

Proof Using $\log p_\theta(x) = \langle \theta, \phi(x) \rangle - A(\theta)$ we obtain $H(P_\theta) = -\mathbb{E}_\theta[\langle \theta, \phi(X) \rangle - A(\theta)] = A(\theta) - \langle \mu(P_\theta), \theta \rangle$, where as usual $\mu(P) = \mathbb{E}_P[\phi(X)]$. As θ and $\mu(P_\theta)$ have the duality relationship $\nabla A(\theta) = \mu(P_\theta)$, we obtain $A(\theta) - \langle \mu(P_\theta), \theta \rangle = -\Omega(\mu(P_\theta))$ as desired. \square

The *maximum entropy principal*, which Jaynes [119] first elucidated in the 1950s, originates in statistical mechanics, where Jaynes showed that (in a sense) entropy in statistical mechanics and information theory were equivalent. The maximum entropy principle is this: given some constraints (prior information) about a distribution P , we consider all probability distributions satisfying said constraints. Then to encode our prior information while being as “objective” or “agnostic” as possible (essentially being as uncertain as possible), we should choose the distribution P satisfying the constraints to maximize the Shannon entropy. This principal naturally gives rise to exponential family models, and (as we revisit later) allows connections to Bayesian and minimax procedures. One caveat throughout is that the base measure ν is *essential* to all our derivations: it radically effects the distributions P we consider.

With all this said, suppose (without making any exponential family assumptions yet) we are given $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$ and a mean vector $\mu \in \mathbb{R}^d$, and we wish to solve

$$\text{maximize } H(P) \quad \text{subject to } \mathbb{E}_P[\phi(X)] = \mu \tag{14.4.2}$$

over all distributions $P \in \mathcal{P}$, the collection of distributions having densities with respect to the base measure ν , that is, $P \ll \nu$. Rewriting problem (14.4.2), we see that it is equivalent to

$$\begin{aligned} & \text{maximize} && -\int p(x) \log p(x) d\nu(x) \\ & \text{subject to} && \int p(x) \phi(x) d\nu(x) = \mu, \quad p(x) \geq 0 \text{ for } x \in \mathcal{X}, \quad \int p(x) d\nu(x) = 1. \end{aligned}$$

Let

$$\mathcal{P}_\mu^{\text{lin}} := \{P \ll \nu \mid \mathbb{E}_P[\phi(X)] = \mu\}$$

be distributions with densities w.r.t. ν satisfying the expectation (linear) constraint $\mathbb{E}[\phi(X)] = \mu$. We then obtain the following theorem.

Theorem 14.4.7. *For $\theta \in \mathbb{R}^d$, let P_θ have density*

$$p_\theta(x) = \exp(\langle \theta, \phi(x) \rangle - A(\theta)), \quad A(\theta) = \log \int \exp(\langle \theta, \phi(x) \rangle) d\nu(x),$$

with respect to the measure ν . If $\mathbb{E}_{P_\theta}[\phi(X)] = \mu$, then P_θ maximizes $H(P)$ over $\mathcal{P}_\mu^{\text{lin}}$; moreover, the distribution P_θ is unique (though θ need not be).

Proof We first give a heuristic derivation—which is not completely rigorous—and then check to verify that our result is exact. First, we write a Lagrangian for the problem (14.4.2). Introducing Lagrange multipliers $\lambda(x) \geq 0$ for the constraint $p(x) \geq 0$, $\theta_0 \in \mathbb{R}$ for the normalization constraint that $P(\mathcal{X}) = 1$, and $\theta \in \mathbb{R}^d$ for the constraints that $\mathbb{E}_P[\phi(X)] = \mu$, we obtain the following Lagrangian:

$$\begin{aligned} \mathcal{L}(p, \theta, \theta_0, \lambda) = & \int p(x) \log p(x) d\nu(x) + \sum_{i=1}^d \theta_i \left(\mu_i - \int p(x) \phi_i(x) d\nu(x) \right) \\ & + \theta_0 \left(\int p(x) d\nu(x) - 1 \right) - \int \lambda(x) p(x) d\nu(x). \end{aligned}$$

Now, heuristically treating the density $p = [p(x)]_{x \in \mathcal{X}}$ as a finite-dimensional vector (in the case that \mathcal{X} is finite, this is completely rigorous), we take derivatives and obtain

$$\frac{\partial}{\partial p(x)} \mathcal{L}(p, \theta, \theta_0, \lambda) = 1 + \log p(x) - \sum_{i=1}^d \theta_i \phi_i(x) + \theta_0 - \lambda(x) = 1 + \log p(x) - \langle \theta, \phi(x) \rangle + \theta_0 - \lambda(x).$$

To find the minimizing p for the Lagrangian (the function is convex in p), we set this equal to zero to find that

$$p(x) = \exp(\langle \theta, \phi(x) \rangle - 1 - \theta_0 - \lambda(x)).$$

Now, we note that with this setting, we always have $p(x) > 0$, so that the constraint $p(x) \geq 0$ is unnecessary and (by complementary slackness) we have $\lambda(x) = 0$. In particular, by taking $\theta_0 = -1 + A(\theta) = -1 + \log \int \exp(\langle \theta, \phi(x) \rangle) d\nu(x)$, we have that (according to our heuristic derivation) the optimal density p should have the form

$$p_\theta(x) = \exp(\langle \theta, \phi(x) \rangle - A(\theta)).$$

So we see the form of distribution we would like to have.

Consider any distribution $P \in \mathcal{P}_\mu^{\text{lin}}$, and assume that we have some θ satisfying $\mathbb{E}_{P_\theta}[\phi(X)] = \mu$. In this case, we may expand the entropy $H(P)$ as

$$\begin{aligned} H(P) &= - \int p \log p d\nu = - \int p \log \frac{p}{p_\theta} d\nu - \int p \log p_\theta d\nu \\ &= -D_{\text{kl}}(P \| P_\theta) - \int p(x) [\langle \theta, \phi(x) \rangle - A(\theta)] d\nu(x) \\ &\stackrel{(\star)}{=} -D_{\text{kl}}(P \| P_\theta) - \int p_\theta(x) [\langle \theta, \phi(x) \rangle - A(\theta)] d\nu(x) \\ &= -D_{\text{kl}}(P \| P_\theta) + H(P_\theta), \end{aligned}$$

where in the step (\star) we have used the fact that $\int p(x)\phi(x)d\nu(x) = \int p_\theta(x)\phi(x)d\nu(x) = \mu$. As $D_{\text{kl}}(P\|P_\theta) > 0$ unless $P = P_\theta$, we have shown that P_θ is the unique distribution maximizing the entropy, as desired. \square

We obtain the following immediate corollary, which shows the direct connection between maximum entropy and minimizing expected logarithmic loss.

Corollary 14.4.8. *Let $\{P_\theta\}$ be the exponential family with densities $p_\theta(x) = \exp(\langle \theta, \phi(x) \rangle - A(\theta))$ with respect to ν . For any $\mu \in \mathcal{M}$, if there exists θ satisfying $\mathbb{E}_{P_\theta}[\phi(X)] = \mu$, then P_θ solves*

$$\underset{p}{\text{minimize}} \mathbb{E}_P[-\log p(x)]$$

over all densities p satisfying $\int \phi(x)p(x)d\nu(x) = \mu$.

So if we consider minimizing the negative log loss (which is strictly proper) but wish to guarantee that the predictive distribution satisfies $\mathbb{E}_P[\phi(X)] = \mu$, then the exponential family model is the unique minimizer.

We give three examples of maximum entropy, showing how the choice of the base measure ν effects the resulting maximum entropy distribution. For all three, we assume that the space $\mathcal{X} = \mathbb{R}$ is the real line. We consider maximizing the entropy over all distributions P satisfying

$$\mathbb{E}_P[X^2] = 1.$$

Example 14.4.9: Assume that the base measure ν is counting measure on the support $\{-1, 1\}$, so that $\nu(\{-1\}) = \nu(\{1\}) = 1$. Then the maximum entropy distribution is given by $P(X = x) = \frac{1}{2}$ for $x \in \{-1, 1\}$. \diamond

Example 14.4.10: Assume that the base measure ν is Lebesgue measure on $\mathcal{X} = \mathbb{R}$, so that $\nu([a, b]) = b - a$ for $b \geq a$. Then by Theorem 14.4.7, we have that the maximum entropy distribution has the form $p_\theta(x) \propto \exp(-\theta x^2)$; recognizing the normal, we see that the optimal distribution is simply $\mathcal{N}(0, 1)$. \diamond

Example 14.4.11: Assume that the base measure ν is counting measure on the integers $\mathbb{Z} = \{\dots, -2, -1, 0, 1, \dots\}$. Then Theorem 14.4.7 shows that the optimal distribution is a discrete version of the normal: we have $p_\theta(x) \propto \exp(-\theta x^2)$ for $x \in \mathbb{Z}$. That is, we choose $\theta > 0$ so that the distribution $p_\theta(x) = \exp(-\theta x^2) / \sum_{j=-\infty}^{\infty} \exp(-\theta j^2)$ has variance 1. \diamond

We remark in passing that in some cases, it is interesting to instead consider *inequality* rather than equality constraints in the linear constraints defining the family \mathcal{P}^{lin} . Exercises 14.10 and 14.11 explore these ideas.

Lastly, we consider the empirical variant of minimizing the log loss, equivalently, of maximum likelihood, where we maximize the likelihood of a given sample X_1, \dots, X_n . Consider the sample-based maximum likelihood problem of solving

$$\underset{\theta}{\text{maximize}} \prod_{i=1}^n p_\theta(X_i) \equiv \underset{\theta}{\text{minimize}} -\frac{1}{n} \sum_{i=1}^n \log p_\theta(X_i), \quad (14.4.3)$$

for the exponential family model $p_\theta(x) = \exp(\langle \theta, \phi(x) \rangle - A(\theta))$. We have the following result.

Proposition 14.4.12. *Let $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \phi(X_i)$. Then any θ solving $\mathbb{E}_{P_\theta}[\phi(X)] = \hat{\mu}_n$ is a maximum likelihood solution, which exists if and only if $\hat{\mu}_n \in \text{relint } \mathcal{M}$. If the sample is drawn $X_i \stackrel{\text{iid}}{\sim} P$ where $P \ll \nu$ and $\mu(P) \in \text{relint } \mathcal{M}$, then with probability 1, $\hat{\mu}_n \in \text{relint } \mathcal{M}$ eventually.*

Proof Define the empirical negative log likelihood

$$\hat{L}_n(\theta) := -\frac{1}{n} \sum_{i=1}^n \log p_\theta(X_i) = -\langle \hat{\mu}_n, \theta \rangle + A(\theta),$$

which is convex. Taking derivatives and using that $\Theta = \text{dom } A$ is open, the parameter θ is a minimizer if and only if $\nabla \hat{L}_n(\theta) = \hat{\mu}_n - \nabla A(\theta) = 0$ if and only if $\nabla A(\theta) = \hat{\mu}_n$. Apply Proposition 14.4.1.

For the final statement, note that $\hat{\mu} \in \text{aff}(\mathcal{M})$ with probability 1. Then because $\mu(P) \in \text{relint } \mathcal{M}$ and $\hat{\mu}_n \rightarrow \mu(P)$ with probability 1, we see that for any $\epsilon > 0$ there is some (random, but finite) N such that $n \geq N$ implies $\|\hat{\mu}_n - \mu(P)\| \leq \epsilon$ and $\hat{\mu}_n \in \text{aff}(\mathcal{M})$, so that $\hat{\mu}_n \in \text{relint } \mathcal{M}$. \square

As a consequence of the result, we have the following rough equivalences tying together the preceding material. In short, maximum entropy subject to (linear) empirical moment constraints (Theorem 14.4.7) is equivalent to maximum likelihood estimation in exponential families (Proposition 14.4.12), and these are all equivalent to minimizing the (surrogate) log loss $\mathbb{E}[\varphi(\theta, X)]$.

14.5 Technical and deferred proofs

14.5.1 Finalizing the proof of Theorem 14.2.15

The issue remaining in the proof of Theorem 14.2.15 occurs when $\ell(\mu, y_i) = +\infty$ for some i . In this case, we necessarily have $p_i = 0$ for all $p \in \Delta_m$ satisfying $\mathbb{E}_p[Y] = \mu$; define the set of infinite loss indices $\mathcal{I}(\mu) := \{i \mid L(\mu, y_i) = +\infty\}$, which is evidently in the set $\{i \mid p_i = 0 \text{ whenever } Ap = 0\}$. Because of this containment, the vectors $\{y_i\}_{i \in \mathcal{I}(\mu)}$ are independent and independent of $\{y_i\}_{i \notin \mathcal{I}(\mu)}$. In particular, there exists $\Delta \in \mathbb{R}^k$ such that $y_i^T \Delta = 0$ for all $i \notin \mathcal{I}(\mu)$ but for which $y_i^T \Delta > 0$ for each $i \in \mathcal{I}(\mu)$. Working on the subspace $\{p \in \Delta_m \mid p_i = 0, i \in \mathcal{I}(\mu)\}$, we can perform precisely the same derivation except that $G(\mu) = \{s \in \mathbb{R}^k \mid y_i^T s = -\ell(\mu, y_i) \text{ for } i \notin \mathcal{I}(\mu)\}$ is non-empty. Then we have

$$\Omega(\mu') = -\mathbb{E}_{p^*(\mu')}[\ell(\mu', Y)] \stackrel{(i)}{\geq} -\mathbb{E}_{p^*(\mu')}[\ell(\mu, Y)] = -\mathbb{E}_{p^*(\mu)}[\ell(\mu, Y)] + \sum_{i=1}^m \ell(\mu, y_i)(p_i^*(\mu) - p_i^*(\mu')),$$

where inequality (i) follows because ℓ is proper. We then have

$$\begin{aligned} \sum_{i=1}^m \ell(\mu, y_i)(p_i^*(\mu) - p_i^*(\mu')) &\stackrel{(ii)}{=} \sum_{i \notin \mathcal{I}(\mu)} \ell(\mu, y_i)(p_i^*(\mu) - p_i^*(\mu')) - \sum_{i \in \mathcal{I}(\mu)} \ell(\mu, y_i)p_i^*(\mu') \\ &= \sum_{i \notin \mathcal{I}(\mu)} s^T y_i(p_i^*(\mu) - p_i^*(\mu')) - \sum_{i \in \mathcal{I}(\mu)} \ell(\mu, y_i)p_i^*(\mu') \end{aligned}$$

for any $s \in G(\mu)$, where equality (ii) follows because $p_i^*(\mu) = 0$ for $i \in \mathcal{I}(\mu)$. As we allow extended reals, replace s with $s_\infty = \lim_{t \rightarrow \infty} (s + t\Delta)$, which satisfies $\langle s_\infty, y_i \rangle = \infty = \ell(\mu, y_i)$ for $i \in \mathcal{I}(\mu)$, and we finally obtain

$$\Omega(\mu') \geq \Omega(\mu) + \sum_{i=1}^m s_\infty^T y_i(p_i^*(\mu) - p_i^*(\mu')) = \Omega(\mu) + \langle s_\infty, \mu - \mu' \rangle.$$

The equality of the loss is as before.

14.5.2 Proof of Proposition 14.4.1

We first give the proof in the case that $\{P_\theta\}$ is a minimal exponential family, meaning that $\langle u, \phi(x) \rangle$ is non-constant in x for each $u \neq 0$, addressing the non-minimal case at the end. Then A is strictly convex (Proposition 3.2.3). As part of this proof, we will show that \mathcal{M}° is indeed open in this case. We show both inclusions $\mathcal{M}^\circ \subset \nabla A(\Theta)$ and that $\nabla A(\Theta) \subset \mathcal{M}^\circ$.

Showing that $\nabla A(\Theta) \subset \mathcal{M}^\circ$. Fix $\theta_0 \in \Theta$, and let $\mu = \nabla A(\theta_0)$. We must show that there exists $\epsilon > 0$ such that for all $\|u\| \leq \epsilon$, the point $\mu + u \in \mathcal{M}$. Let $\theta_u = \operatorname{argmin}_\theta \{A(\theta) - \langle \mu + u, \theta \rangle\}$ whenever the minimizer exists, where evidently θ_0 does exist because $\mu = \nabla A(\theta_0)$. Note that the strict convexity of A guarantees θ_u is unique if it exists. But now, we may use the convex analytic fact (Proposition C.1.12 in Appendix C.1.2) that $u \mapsto \theta_u$ is continuous in u in a neighborhood of 0. These minimizers necessarily satisfy $\nabla A(\theta_u) = \mu + u$, that is, $\mathbb{E}_{\theta_u}[\phi(X)] = \mu + u \in \mathcal{M}$.

Showing that $\mathcal{M}^\circ \subset \nabla A(\Theta)$. Let $\mu \in \mathcal{M}^\circ$, so that there exists an $\epsilon > 0$ such that $\mu + \epsilon \mathbb{B} \subset \mathcal{M}^\circ$. It is enough to show that $A(\theta) - \langle \mu, \theta \rangle$ is coercive in θ , as then there necessarily exists a (unique) minimizer $\theta(\mu)$ of $A(\theta) - \langle \mu, \theta \rangle$, and this minimizer satisfies $\nabla A(\theta(\mu)) = \mu$, so that $\mu \in \nabla A(\Theta)$. For this, it is sufficient to show that for any non-zero vector v the *recession function* of the tilted version $f(\theta) := A(\theta) - \langle \mu, \theta \rangle$ of A ,

$$f'_\infty(v) := \lim_{t \rightarrow \infty} \frac{A(\theta + tv) - \langle \mu, \theta + tv \rangle - (A(\theta) - \langle \mu, \theta \rangle)}{t}$$

where $\theta \in \Theta$ is otherwise arbitrary, satisfies $f'_\infty(v) > 0$ for all $v \neq 0$, which guarantees that $A(\cdot) - \langle \mu, \cdot \rangle$ has a minimizer. (See Proposition C.3.5 and Corollary C.3.7 in Appendix C.2.1).

To that end, for vectors $v \in \mathbb{R}^d$, define the essential supremum of $\phi(x)$ in the direction v by

$$\nu^*(\phi, v) := \operatorname{ess\,sup}_x \langle \phi(x), v \rangle = \inf_t \{t \in \mathbb{R} \mid \nu(\{x \in \mathcal{X} \mid \langle v, \phi(x) \rangle \geq t\}) = 0\}.$$

Now as $\mu \in \mathcal{M}^\circ$, for any vector $v \neq 0$ we have $\langle v, \mu \rangle < \nu^*(\phi, v)$. Let $\epsilon > 0$ satisfy $\langle v, \mu \rangle < \nu^*(\phi, v) - \epsilon$ be otherwise arbitrary, fix $\theta \in \Theta$, and let $\mathcal{X}_\epsilon = \{x \mid \langle v, \phi(x) \rangle \geq \nu^*(\phi, v) - \epsilon\}$, which satisfies $\nu(\mathcal{X}_\epsilon) > 0$. Then

$$\begin{aligned} A(\theta + tv) - \langle \mu, \theta + tv \rangle &= \log \int \exp(\langle \phi(x), \theta + tv \rangle) d\nu(x) - \langle \mu, \theta + tv \rangle \\ &\geq \log \int_{\mathcal{X}_\epsilon} \exp(\langle \phi(x), \theta \rangle) e^{t(\nu^* - \epsilon)} d\nu(x) - \langle \mu, \theta \rangle - t\langle \mu, v \rangle \\ &= t(\nu^*(\phi, v) - \epsilon) + \log \nu(\mathcal{X}_\epsilon) - t\langle \mu, v \rangle + \log \int_{\mathcal{X}_\epsilon} e^{\langle \phi(x), \theta \rangle} d\nu(x) - \langle \mu, \theta \rangle. \end{aligned}$$

If $\nu(\mathcal{X}_\epsilon) = +\infty$, then $A(\theta + tv) = +\infty$ and so $A'_\infty(v) > 0$ certainly. If $\nu(\mathcal{X}_\epsilon) < \infty$, then note that $\nu^*(\phi, v) - \epsilon - \langle \mu, v \rangle > 0$, and so

$$A(\theta + tv) - \langle \mu, \theta + tv \rangle \geq t(\nu^*(\phi, v) - \epsilon - \langle \mu, v \rangle) - \log \nu(\mathcal{X}_\epsilon) + \log \int_{\mathcal{X}_\epsilon} e^{\langle \phi(x), \theta \rangle} d\nu(x) - \langle \mu, \theta \rangle$$

and thus

$$\frac{A(\theta + tv) - \langle \mu, \theta + tv \rangle - (A(\theta) - \langle \mu, \theta \rangle)}{t} \geq \nu^*(\phi, v) - \epsilon - \langle \mu, v \rangle + o(1) \quad (14.5.1)$$

as $t \rightarrow \infty$.

Extending to the non-minimal case. If the exponential family is not minimal, there exists a unit vector u and constant c such that $\langle u, \phi(x) \rangle = c$ for ν -almost all x . Let $U \in \mathbb{R}^{d \times k}$ be an orthonormal basis for all such vectors, where k is the dimension of this collection. Then there exists a vector $c \in \mathbb{R}^k$ such that $c = U^\top \phi(x)$ for ν -almost all x , and we see that $A(\theta + Uv) = A(\theta) + \langle c, v \rangle$ as $\langle \theta + Uv, \phi(x) \rangle = \langle \theta, \phi(x) \rangle + \langle c, v \rangle$ for ν -almost all x . We show both inclusions as above. Let $U_\perp \in \mathbb{R}^{d \times d-k}$ be an orthonormal basis for the orthogonal subspace to U , so that $U^\top U = I_k$ and $U_\perp^\top U_\perp = I_{d-k}$, and for any $\mu \in \mathcal{M}$, we have $\text{aff}(\mathcal{M}) = \mu + \text{span}(U_\perp)$.

Showing that $\nabla A(\Theta) \subset \mathcal{M}^\circ$. Fix $\theta_0 \in \Theta$ and let $\mu = \nabla A(\theta_0)$. We must show that there exists $\epsilon > 0$ such that for all $u \in \text{span}(U_\perp)$ satisfying $\|u\| \leq \epsilon$, the point $\mu + u \in \mathcal{M}$. To that end, note that for any vectors $v \in \mathbb{R}^{d-k}$ and $w \in \mathbb{R}^k$, we have

$$A(\theta_0 + U_\perp v + Uw) - \langle \mu + u, U_\perp v + Uw \rangle = A(\theta_0 + U_\perp v) - \langle \mu + u, U_\perp v \rangle$$

because $U^\top u = 0$ and $U^\top \mu = c$ for each $u \in \text{span}(U_\perp)$ and $\mu \in \mathcal{M}$. The function $g(v) := A(\theta_0 + U_\perp v) - \langle \mu, U_\perp v \rangle$ is strictly convex as $\nabla^2 g(v) = U_\perp^\top \nabla^2 A(\theta_0 + U_\perp v) U_\perp \succ 0$, because we know that $u^\top \phi(x)$ is non-constant for all $u \in \text{span}(U_\perp)$. Define $f(v) = A(\theta_0 + U_\perp v) - \langle \mu, U_\perp v \rangle$. Then applying Proposition C.1.12 as in the minimal representation case, there exists $\epsilon > 0$ such that $v_u = \text{argmin}_v \{f(v) - \langle u, U_\perp v \rangle\}$ exists and is continuous in $u \in \text{span}(U_\perp)$, where by inspection $v_0 = 0$. Then $\theta_u := \theta_0 + U_\perp v_u$ minimizes $A(\theta) - \langle \mu + u, \theta \rangle$, satisfying $\nabla A(\theta_u) = \mu + u$.

Showing that $\mathcal{M}^\circ \subset \nabla A(\Theta)$. We again follow the logic of the minimal representation case. Let $\mu \in \mathcal{M}^\circ = \text{relint } \mathcal{M}$, and recall $\nu^*(\phi, U_\perp v) = \text{ess sup}_x \langle \phi(x), U_\perp v \rangle$. Then there exists $\epsilon > 0$ such that $\mu + u \in \mathcal{M}$ for each $u \in \text{span}(U_\perp)$ with $\|u\| \leq \epsilon$, so that

$$\langle \mu, U_\perp v \rangle < \sup_{\|u\|_2 \leq \epsilon, u \in \text{span}(U_\perp)} \langle \mu + u, U_\perp v \rangle \leq \nu^*(\phi, U_\perp v).$$

Define $g(v) = A(\theta + U_\perp v) - \langle \mu, U_\perp v \rangle$. Then because $A(\theta + Uw + U_\perp v) - \langle \mu, U_\perp v - Uw \rangle = A(\theta + U_\perp v) - \langle \mu, U_\perp v \rangle$ for all $w \in \mathbb{R}^k, v \in \mathbb{R}^{d-k}$, it is enough to show that $g'_\infty(v) > 0$ for all $v \neq 0$. Following the same argument, *mutatis mutandis*, as that leading to inequality (14.5.1) yields that $g'_\infty(v) > 0$ for all $v \neq 0$. That is, $v \mapsto A(\theta + U_\perp v) - \langle \mu, U_\perp v \rangle$ has a minimizer $v(\mu)$ (Corollary C.3.7), which is unique by the strict convexity of $v \mapsto A(\theta + U_\perp v)$, and which necessarily satisfies $U_\perp^\top \nabla A(\theta + U_\perp v(\mu)) = U_\perp^\top \mu$. As $U^\top \nabla A(\theta) = c$ for all θ and $U^\top \mu = c$ for all $\mu \in \mathcal{M}$, this shows that there exists $\theta(\mu)$ such that $\nabla A(\theta(\mu)) = \mu$ as desired. Moreover, fixing an arbitrary θ and letting $v(\mu)$ be the unique minimizer of $A(\theta + U_\perp v) - \langle \mu, U_\perp v \rangle$, the set of all minimizers

$$\Theta^*(\mu) = \text{argmin}_\theta \{A(\theta) - \langle \mu, \theta \rangle\} = \left\{ \theta + U_\perp v(\mu) + Uw \mid w \in \mathbb{R}^k \right\}.$$

This gives Proposition 14.4.1.

14.5.3 Proof of Proposition 14.4.3

For part (i), because $\Theta = \text{dom } A \subset \mathbb{R}^d$ is open and A is \mathcal{C}^∞ on its domain, A is necessarily a closed convex function and so $A^{**}(\theta) = A(\theta)$ for all $\theta \in \mathbb{R}^d$. (See Theorem C.2.1.) For part (ii), note that if $\mu \in \mathcal{M}^\circ$, there exists $\theta(\mu) \in \Theta$ such that $\nabla A(\theta(\mu)) = \mu$ by Proposition 14.4.1. This $\theta(\mu)$ maximizes $\langle \theta, \mu \rangle - A(\theta)$ over all θ , and so $\Omega(\mu) = \langle \theta(\mu), \mu \rangle - A(\theta(\mu)) < \infty$. By Corollary C.2.5 in Appendix C.2.1 and Proposition 14.4.1, $\text{dom } \partial\Omega = \mathcal{M}^\circ$, and as Ω is subdifferentiable on the relative interior of its domain, we have $\text{dom } \Omega \subset \text{cl } \mathcal{M}^\circ = \text{cl } \mathcal{M}$. As Ω is closed convex, any point μ outside its domain necessarily satisfies $\Omega(\mu) = +\infty$.

Finally, for part (iii), we note that the function $g(t) = \Omega(t\mu_0 + (1-t)\mu)$ is a one-dimensional closed convex function. One-dimensional closed convex functions are continuous on their domains (Observation B.3.6 in Appendix B.3.2), and so g is necessarily continuous. Thus $\lim_{t \downarrow 0} g(t) = g(0)$. The existence of θ_t follows from Proposition 14.4.1.

Bibliography

JCD Comment: Need to do a lot here!

Gneiting and Raftery [103]

14.6 Exercises

Exercise 14.1 (Strict propriety of the log loss): Let $\Delta_k = \{p \in \mathbb{R}_+^k \mid \mathbf{1}^T p = 1\}$ be the probability simplex. Show that if $\ell(q, y) = -\log q_y$ and $\mathbb{P}(Y = y) = p_y$, then

$$\operatorname{argmin}_{q \in \Delta_k} \mathbb{E}[\ell(q, Y)] = p,$$

where we treat $0 \log 0$ as 0 (which is the natural limit of $t \log t$ as $t \downarrow 0$).

Exercise 14.2 (Uniqueness of generalized entropies): Here we give an alternative perspective on the generalized entropies associated with losses, showing when they are unique. For a concave function $f : \Delta_k \rightarrow \mathbb{R}$, define the perspective-type transform $f_{\text{per}}(p) = \langle \mathbf{1}, p \rangle f(p / \langle \mathbf{1}, p \rangle)$, where $f_{\text{per}}(0) = 0$, and which gives $f_{\text{per}} : \mathbb{R}_+^k \rightarrow \mathbb{R}$.

- (a) Let $\ell : \Delta_k \rightarrow \mathbb{R}$ be strictly proper and let Y have p.m.f. p . Show that $H(p) = \inf_{q \in \Delta_k} \mathbb{E}[\ell(q, Y)]$ is strictly concave, and that H_{per} is strictly concave and continuously differentiable on \mathbb{R}_{++}^k .
- (b) Show the converse that if $H : \Delta_k \rightarrow \mathbb{R}$ is strictly concave and its perspective H_{per} is differentiable on \mathbb{R}_{++}^k , then there exists a proper scoring loss ℓ satisfying

$$H(p) = \inf_{q \in \Delta_k} \mathbb{E}_p[\ell(q, Y)]$$

and that $\ell(q, y) = \nabla_y H_{\text{per}}(q)$ for all $q \in \operatorname{dom} \nabla H_{\text{per}}$.

Exercise 14.3: Give the details in the computations for Example 14.3.4.

Exercise 14.4: Let $y \in \{0, 1\}$ and take the regularization function $\Omega(p) = -\log p - \log(1 - p)$.

- (a) Verify that the entropy is of Legendre type (Definition 14.4).
- (b) Give the associated loss ℓ and surrogate loss φ in the sense of Section 14.3.
- (c) Plot the surrogate $\varphi(s, y) + \log 8$ and the logistic regression surrogate $\log(1 + e^s) - sy$ for $y \in \{0, 1\}$, each as function of s . (The shift by $\log 8$ guarantees the losses coincide at $s = 0$.)
- (d) Give $\operatorname{pred}_\Omega(s)$ for $s \in \mathbb{R}$, verifying that $\operatorname{pred}_\Omega(s) \in [0, 1]$.

Exercise 14.5: For $\Omega(p) = -\log p - \log(1 - p)$ as in Exercise 14.4, show that Ω is *self-concordant*, meaning that $\Omega'''(p) \leq 2(\Omega''(p))^{3/2}$ for all $p \in (0, 1)$. (Such functions are important in optimization; the conjugate Ω^* is then also guaranteed to be self-concordant.)

Exercise 14.6 (Surrogates for regression): Define $\Omega(\mu) = \frac{1}{4}\mu^4$.

- (a) Give the conjugate $\Omega^*(s)$ to Ω .
- (b) Show directly that the surrogate loss $\varphi(s, y) = \Omega^*(s) - sy$ satisfies that if $\hat{s} = \operatorname{argmin}_s \mathbb{E}[\varphi(s, Y)]$, then $\operatorname{pred}_\Omega(\hat{s}) = \mathbb{E}[Y]$.

Exercise 14.7: Let P be a predicted distribution and for $\alpha \in [0, \frac{1}{2}]$, define the lower and upper quantiles $l_\alpha = \operatorname{Quant}_\alpha(P)$ and $u_\alpha = \operatorname{Quant}_{1-\alpha}(P)$. Given these quantiles, for a finite set $\mathcal{A} \subset [0, \frac{1}{2}]$, define the weighted interval loss

$$W(P, y) := \sum_{\alpha \in \mathcal{A}} [\alpha(u_\alpha - l_\alpha) + \operatorname{dist}(y, [l_\alpha, u_\alpha])],$$

which penalizes P using both the size $(u_\alpha - l_\alpha)$ of the quantile intervals and the distance of the outcome y from the predicted quantiles. Define the symmetrized set $\mathcal{A}_s = \mathcal{A} \cup \{1 - \alpha \mid \alpha \in \mathcal{A}\}$. Show that

$$W(P, y) = \ell_{\operatorname{quant}, \mathcal{A}_s}(P, y),$$

where $\ell_{\operatorname{quant}}$ is the quantile loss (14.2.4).

Exercise 14.8: We explore a particularization of the results in Section 14.4. Let $Y \sim \operatorname{Poi}(e^\theta)$, so that Y has p.m.f. $p_\theta(y) = \exp(\theta y - e^\theta)/y!$ for $y \in \mathbb{N}$. Let $A(\theta) = e^\theta$ be the log-partition function. Define the “surrogate” loss $\varphi(\theta, y) = -\log p_\theta(y)$.

- (a) Give the associated negative generalized entropy $\Omega(\mu)$ for $\mu \in (0, \infty)$.
- (b) Give the associated loss $\ell(\mu, y)$ in the proper representation of Theorem 14.2.15. Directly verify that it is strictly proper, in that $\operatorname{argmin}_\mu \mathbb{E}[\ell(\mu, Y)] = \mathbb{E}[Y]$ for any Y supported on \mathbb{R}_+ .

Exercise 14.9: We explore a particularization of Example 14.4. Let $X \sim \mathcal{N}(0, \Sigma)$ for a covariance $\Sigma \succ 0$, and let $K = \Sigma^{-1}$ be the associated precision matrix. Then X has density $p_K(x) = \exp(-\frac{1}{2}\langle xx^T, K \rangle + \frac{1}{2} \log \det(K))$ with respect to (a scaled) Lebesgue measure, and log partition $A(K) = -\frac{1}{2} \log \det(K)$, which has domain the positive definite matrices $K \succ 0$ (and is $+\infty$ elsewhere).

- (a) Give the associated negative generalized entropy $\Omega(M)$ for symmetric matrices M . Specify the domain of Ω .
- (b) Give the associated loss $\ell(M, x)$ in the proper representation of Theorem 14.2.15. Directly verify that it is strictly proper, in that if the second moment matrix $C := \mathbb{E}[XX^T]$ of X satisfies $C \succ 0$, then $\operatorname{argmin}_M \mathbb{E}[\ell(M, X)] = C$.

Exercise 14.10: In this extended exercise, we generalize Theorem 14.4.7 to apply to general (finite-dimensional) convex cone constraints. A set \mathcal{C} is a *convex cone* if for any two points $x, y \in \mathcal{C}$, we have $\lambda x + (1 - \lambda)y \in \mathcal{C}$ for all $\lambda \in [0, 1]$, and \mathcal{C} is closed under positive scaling: $x \in \mathcal{C}$ implies that $tx \in \mathcal{C}$ for all $t \geq 0$. The following are standard examples (the positive orthant and the semi-definite cone):

- i. *The orthant.* Take $\mathcal{C} = \mathbb{R}_+^d = \{x \in \mathbb{R}^d : x_j \geq 0, j = 1, \dots, d\}$. Then clearly \mathcal{C} is convex and closed under positive scaling.

- ii. *The semidefinite cone.* Take $\mathcal{C} = \{X \in \mathbb{R}^{d \times d} : X = X^\top, X \succeq 0\}$, where a matrix $X \succeq 0$ means that $a^\top X a \geq 0$ for all vectors a . Then \mathcal{C} is convex and closed under positive scaling as well.

Given a convex cone \mathcal{C} , we associate a cone ordering \succeq with the cone and say that for two elements $x, y \in \mathcal{C}$, we have $x \succeq y$ if $x - y \in \mathcal{C}$, that is, $x - y \in \mathcal{C}$. In the orthant case, this simply means that x is component-wise larger than y . For a given inner product $\langle \cdot, \cdot \rangle$, define the dual cone

$$\mathcal{C}^* := \{y : \langle y, x \rangle \geq 0 \text{ for all } x \in \mathcal{C}\}.$$

For the standard (Euclidean) inner product, the positive orthant is thus self-dual, and similarly the semidefinite cone is also self-dual. For a vector y , we write $y \succeq_* 0$ if $y \in \mathcal{C}^*$ is in the dual cone. With this setup, consider the following linearly constrained maximum entropy problem, where the cone ordering \preceq derives from a cone \mathcal{C} :

$$\text{maximize } H(P) \quad \text{subject to } \mathbb{E}_P[\phi(X)] = \mu, \quad \mathbb{E}_P[\psi(X)] \preceq \beta, \quad (14.6.1)$$

where the base measure ν is implicit. Let $\mathcal{P}_{\mu, \beta}^{\text{lin}}$ be the collection of distributions $P \ll \nu$ satisfying the constraints in problem (14.6.1).

Prove the following theorem:

Theorem 14.6.1. *For $\theta \in \mathbb{R}^d$ and $K \in \mathcal{C}^*$, the dual cone to \mathcal{C} , let $P_{\theta, K}$ have density*

$$p_{\theta, K}(x) = \exp(\langle \theta, \phi(x) \rangle - \langle K, \psi(x) \rangle - A(\theta, K)), \quad A(\theta, K) = \log \int \exp(\langle \theta, \phi(x) \rangle - \langle K, \psi(x) \rangle) d\nu(x),$$

with respect to the measure ν . If

$$\mathbb{E}_{P_{\theta, K}}[\phi(X)] = \mu \quad \text{and} \quad \mathbb{E}_{P_{\theta, K}}[\psi(X)] = \beta,$$

then $P_{\theta, K}$ maximizes $H(P)$ over $\mathcal{P}_{\mu, \beta}^{\text{lin}}$. Moreover, the distribution $P_{\theta, K}$ is unique.

Exercise 14.11 (An application of Theorem 14.6.1): Let the cone \mathcal{C} be the positive semidefinite cone in $\mathbb{R}^{d \times d}$, ν be the Lebesgue measure $d\nu(x) = dx$ and define $\psi(x) = \frac{1}{2}xx^\top \in \mathbb{R}^{d \times d}$. Let $\Sigma \succ 0$. Give the density solving

$$\text{maximize } - \int p(x) \log p(x) dx \quad \text{subject to } \mathbb{E}_P[XX^\top] \preceq \Sigma.$$

Exercise 14.12: Prove that the log determinant function is concave over the positive semidefinite matrices. That is, show that for $X, Y \in \mathbb{R}^{d \times d}$ satisfying $X \succeq 0$ and $Y \succeq 0$, we have

$$\log \det(\lambda X + (1 - \lambda)Y) \geq \lambda \log \det(X) + (1 - \lambda) \log \det(Y)$$

for any $\lambda \in [0, 1]$. *Hint:* think about log-partition functions.

Exercise 14.13 (Entropy and log-determinant maximization): Consider the following optimization problem over symmetric positive semidefinite matrices in $\mathbb{R}^{d \times d}$:

$$\text{maximize}_{\Sigma \succeq 0} \log \det(\Sigma) \quad \text{subject to } \Sigma_{ij} = \sigma_{ij}$$

where σ_{ij} are specified only for indices $i, j \in S$ (but we know that $\sigma_{ij} = \sigma_{ji}$ and $(i, i) \in S$ for all i). Let Σ^* denote the solution to this problem, assuming there is a positive definite matrix Σ satisfying

$\Sigma_{ij} = \sigma_{ij}$ for $i, j \in S$. Show that for each unobserved pair $(i, j) \notin S$, the (i, j) entry $[\Sigma^{*-1}]_{ij}$ of the inverse Σ^{*-1} is 0. *Hint:* The distribution maximizing the entropy $H(X) = -\int p(x) \log p(x) dx$ subject to $\mathbb{E}[X_i X_j] = \sigma_{ij}$ has Gaussian density of the form $p(x) = \exp(\sum_{(i,j) \in S} \lambda_{ij} x_i x_j - \Lambda_0)$.

Exercise 14.14: **JCD Comment:** Finish this.

Equivalence of integrated quantile losses and continuous ranked probability score.

Chapter 15

Calibration and Proper Losses

In Chapter 14, we encountered *proper losses*, in which we assume we predict probability distributions on outcomes Y . In typical problems, we wish to predict things about Y from a given set of covariates or inputs X , and in focusing exclusively on the losses ℓ themselves, we implicitly assume that we can model $Y | X$ basically perfectly. Here, we move away from this focus exclusively on the loss itself to incorporate discussion of predictions, where we seek a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ (or some other output space) that yields the most accurate predictions.

In this chapter, we adopt the view of Section 14.2.3, where the target $Y \subset \mathbb{R}^k$ is vector-valued, and we wish to predict its expectation $\mathbb{E}[Y | X]$ as accurately as possible. For binary prediction, we have $Y \in \{0, 1\}$, so that $\mathbb{E}[Y | X] = \mathbb{P}(Y = 1 | X)$; in the case of multiclass prediction problems, it is easy to represent Y as an element of the k standard basis vectors $\{e_1, \dots, e_k\} \subset \mathbb{R}^k$, so that $p = \mathbb{E}[Y | X]$ is simply the p.m.f. of Y given X with entries $p_y = \mathbb{P}(Y = y | X)$. We focus here, therefore, on choosing functions to minimize the risk, or expected population loss,

$$R(f) := \mathbb{E}[\ell(f(X), Y)].$$

When f is chosen from a collection $\mathcal{F} \subset \{\mathcal{X} \rightarrow \mathbb{R}^k\}$ of functions, for example, to guarantee that we can generalize, we do not expect to be able to perfectly minimize the population loss. Accordingly, even though the loss is proper and hence minimized by $f^*(x) = \mathbb{E}[Y | X = x]$, we cannot perfectly model reality, and so it is unrealistic to expect to be able to find f satisfying $f(x) = \mathbb{E}[Y | X = x]$, even approximately, for all x .

We therefore depart from the goal of perfection to address a somewhat simpler criterion: that of calibration. Here, the idea is that a predictor should be accurate on average conditional on its own predictions. Consider again a weather forecasting problem, where $Y_t = 1$ indicates it rains on day t and $Y_t = 0$ indicates no rain, and we wish to predict Y_t based on observable covariates X_t at time t . While we would like a forecaster to have perfect predictions $p_t = \mathbb{E}[Y_t | X_t]$, we instead ask that on days where the forecaster makes a given prediction, it should rain (roughly) with that given frequency. In particular, we seek *calibration*, which is that

$$f(X) = \mathbb{E}[Y | f(X)]. \tag{15.0.1}$$

That is, given that the forecaster makes a prediction with value $p = f(X)$, we should have

$$\mathbb{E}[Y | f(X) = p] = p.$$

While in general it is challenging to achieve this perfect calibration, in this chapter we investigate several variants of the desideratum (15.0.1) that allow for more elegant statistical and information-theoretic approaches, as well as procedures to achieve calibration.

This chapter therefore proceeds as follows. The first goal is to

JCD Comment: Fix notation. Also add a transition here to make clearer why we are doing this and what we are doing.

1. First show what we want to measure.
2. Show how to measure it, specifically using partitioned methods. I think that partitioned ones should be better than non-partitioned approaches, because we can estimate the binned / partitioned calibration error
3. Show a few ways to achieve it (population and finite-sample level).

It is important to note that the literature on calibration is broad, and there are several distinct strands. We take the particular focus that most dovetails with our treatment of proper losses and scoring rules, basing our development around random variables and finite-dimensional probabilities. So, for example, if a logistic regression model (as in Example 3.4.2 or 3.4.3) for image classification assigns a probability of 80% that an image is, say, a dog, then the model is (approximately) calibrated if in the population of all images in the world to which the model assigns probability 80%, (approximately) 80% are dogs. The first direction of research that we essentially do not touch are the following: in the forecasting literature, one often considers predicting the distribution of a (potentially continuous) random variable Y , such as the amount of rainfall; if we predict a cumulative distribution F as in Example 14.2.6, then perfect calibration (15.0.1) becomes that

$$\mathbb{P}(Y \leq u \mid F) = F(u) \quad \text{for all } u \in \mathbb{R}.$$

This is far too stringent a condition to be achievable, so that one relaxes to various forms of marginal or average calibration. See the bibliographic notes for some discussion of the approaches here.

The second strand of research on calibration that, again, we do not address, considers more adversarial and sequential settings, where instead of any probabilistic underpinnings, nature (an adversary) plays a game against the player (or predictor). Philosophically, this approach elegantly does away with the need for probabilities: there is a physical world where whether it rains tomorrow is essentially deterministic, and we use probability as a crutch to model things we cannot measure, so calibration means that of the days on which we predict rain with a chance of 50%, it rains on roughly 50% of those days. In this sequential setting, at times $t = 1, 2, \dots, T$, the player makes a prediction p_t of the outcome, and then nature may choose the outcome Y_t . Without giving the player a bit more leeway, calibration is impossible: say that $Y \in \{0, 1\}$, and nature plays $Y_t = 1$ if $p_t \leq .5$ and $Y_t = 0$ if $p_t > .5$. Then any player is miscalibrated at least by an amount .5. Astoundingly, Foster and Vohra [90] show that if the player is allowed to randomize, then the forecasted probabilities p_t can be made arbitrarily close to the empirical averages of the observed Y_t . While many of the techniques we consider and develop arise from this adversarial setting in the literature, we shall mostly address the scenarios in which Y is indeed random.

15.1 Proper losses and calibration error

When we use a proper loss to measure the error $\ell(f(x), y)$ in making the prediction $f(x)$ for the value y , it turns out we can *always* improve the losses by modifying f to be a calibrated version of itself: calibration is always useful. To make this precise, assume we are making predictions in the convex hull of \mathcal{Y} , that is, that can be represented as $\mathbb{E}[Y]$ for some distribution, so $f : \mathcal{X} \rightarrow \mathcal{M} = \text{Conv}(\mathcal{Y})$.

Then by Theorems 14.2.1 and 14.2.15, there exists a convex function h such that

$$\ell(\mu, y) = -h(\mu) - \langle \nabla h(\mu), y - \mu \rangle \quad (15.1.1)$$

for all $\mu \in \mathcal{M}, y \in \mathcal{Y}$. Recall the Bregman divergence (14.2.2)

$$D_h(u, v) = h(u) - h(v) - \langle \nabla h(v), u - v \rangle,$$

which is nonnegative for all convex h (and strictly positive whenever h is strictly convex and $u \neq v$), and Corollary 14.2.5. Then for any prediction function f , if we condition on the predicted value $S = f(X)$, then

$$\begin{aligned} \mathbb{E}[\ell(S, Y) \mid S] &= \mathbb{E}[\ell(\mathbb{E}[Y \mid S], Y) \mid S] + \mathbb{E}[\ell(S, Y) - \ell(\mathbb{E}[Y \mid S], Y) \mid S] \\ &= \mathbb{E}[\ell(\mathbb{E}[Y \mid S], Y) \mid S] + \mathbb{E}[D_h(\mathbb{E}[Y \mid S], S) \mid S], \end{aligned}$$

where we use the linearity $\mathbb{E}[\ell(s, Y)] = \ell(s, \mathbb{E}[Y])$ for any distribution on Y and fixed $s \in \mathbb{R}^k$ in the second equality. We record this as a theorem.

Theorem 15.1.1. *Let ℓ be a proper loss with representation (15.1.1). Then for any $f : \mathcal{X} \rightarrow \mathbb{R}^k$,*

$$\mathbb{E}[\ell(f(X), Y)] = \mathbb{E}[\ell(\mathbb{E}[Y \mid f(X)], Y)] + \mathbb{E}[D_h(\mathbb{E}[Y \mid f(X)], f(X))].$$

In particular, the predictor $g : \mathbb{R}^k \rightarrow \mathbb{R}^k$ defined by

$$g(s) := \mathbb{E}[Y \mid f(X) = s]$$

is calibrated and satisfies

$$\mathbb{E}[\ell(g \circ f(X), Y)] = \mathbb{E}[\ell(\mathbb{E}[Y \mid f(X)], Y)] \leq \mathbb{E}[\ell(f(X), Y)],$$

and the inequality is strict whenever f is not calibrated and ℓ is strictly proper.

Proof The first statement we have already proved. For the second, note that

$$g(s) = \mathbb{E}[Y \mid f(X) = s]$$

by construction of g , so that $\mathbb{E}[\ell(g \circ f(X), Y)] = \mathbb{E}[\ell(\mathbb{E}[Y \mid f(X)], Y)]$. The inequality and its strictness are immediate because h is strictly convex if and only if ℓ is strictly proper. \square

To interpret this result, it essentially says that if we can post-process f to make it calibrated, then we can only improve its risk, or expected loss, when ℓ is a proper loss. We can give an alternative version of Theorem 15.1.1, where we instead consider the conjugate linkages in Section 14.3, which can be useful when we wish to find f via convex optimization (instead of by directly minimizing a proper loss). To that end, assume that h is a strictly convex function, differentiable on the interior of its domain, satisfying the Legendre conditions (14.3.3), and define the surrogate loss (linked via duality and the negative generalized entropy h to ℓ)

$$\varphi(s, y) = h^*(s) - \langle s, y \rangle = \ell(\text{pred}_h(s), y),$$

where $\ell(\mu, y) = -h(\mu) - \langle \nabla h(\mu), y - \mu \rangle$ and

$$\text{pred}_h(s) = \underset{\mu}{\operatorname{argmin}} \{ -\langle s, \mu \rangle + h(\mu) \} = \nabla h^*(s).$$

Then we have the following decomposition of the population surrogate loss, which follows similarly to Theorem 15.1.1.

Theorem 15.1.2. *Let φ be the surrogate loss defined above. Then for any $f : \mathcal{X} \rightarrow \mathbb{R}^k$, we have*

$$\mathbb{E}[\varphi(f(X), Y)] = \mathbb{E}[\ell(\mathbb{E}[Y | f(X)], Y)] + \mathbb{E}[D_h(\mathbb{E}[Y | f(X)], \text{pred}_h(f(X)))].$$

Proof The key is to rely on the duality relationships inherent in the definition of the surrogate $\varphi(s, y) = h^*(s) - \langle s, y \rangle$. We fix x and work exclusively in the space of the scores (predictions) $s = f(x) \in \mathbb{R}^k$, as

$$\mathbb{E}[\varphi(f(X), Y) | X = x] = \varphi(f(x), \mathbb{E}[Y | X = x])$$

by definition. Let $\mu \in \mathcal{M} = \text{Conv}(\mathcal{Y})$. Then $\varphi(s, \mu) = h^*(s) - \langle s, \mu \rangle$, and

$$\inf_s \varphi(s, \mu) = -\sup_s \{\langle s, \mu \rangle - h^*(s)\} = -h(\mu)$$

because h is (closed) convex. Additionally, if $\mu^*(s) = \nabla h^*(s) = \text{pred}_h(s)$, then the conjugate duality relationships (14.1.4) guarantee $h^*(s) = \langle s, \mu^*(s) \rangle - h(\mu^*(s))$ and $s = \nabla h(\mu^*(s))$. Thus

$$\begin{aligned} \varphi(s, \mu) - \inf_{s'} \varphi(s', \mu) &= h^*(s) - \langle s, \mu \rangle + h(\mu) = h(\mu) - h(\mu^*(s)) - \langle s, \mu - \mu^*(s) \rangle \\ &= h(\mu) - h(\mu^*(s)) - \langle \nabla h(\mu^*), \mu - \mu^*(s) \rangle = D_h(\mu, \mu^*(s)). \end{aligned}$$

Taking the expectation over X and using the shorthand $S = f(X)$, we thus obtain

$$\begin{aligned} \mathbb{E}[\varphi(S, Y)] &= \mathbb{E}[\varphi(S, \mathbb{E}[Y | S])] \\ &= \mathbb{E} \left[\inf_s \varphi(s, \mathbb{E}[Y | S]) \right] + \mathbb{E} [D_h(\mathbb{E}[Y | S], \text{pred}_h(S))]. \end{aligned}$$

Lastly, we use that $\ell(\mu, y) = -h(\mu) - \langle \nabla h(\mu), y - \mu \rangle$ is proper, so $\inf_s \varphi(s, \mu) = -h(\mu) = \ell(\mu, \mu)$, giving the first claim of the theorem. \square

As in Theorem 15.1.1, Theorem 15.1.2 shows that calibrating a predictor f can only improve the surrogate loss associated with h . Any predictor $f : \mathcal{X} \rightarrow \mathbb{R}^k$ has unnecessary error arising from the average divergence of the prediction from being calibrated,

$$\mathbb{E}[D_h(\mathbb{E}[Y | f(X)], \text{pred}_h(f(X)))].$$

In both cases, we see that any proper (or derived proper) loss has a natural decomposition into an error term relating to the typical error in predicting Y from $\mathbb{E}[Y | f(X)]$, which one frequently refers to as *sharpness* of the predictor. Replacing $f(X)$ with the expectation of Y given $f(X)$ (or a particular transformation thereof) does not increase this first term, but improves the second term, which measures the typical error of a prediction from calibration.

Let us consider an example with squared error:

Example 15.1.3 (Squared error and calibration): In the case that $h(p) = \frac{1}{2} \|p\|_2^2$, we have $h^* = h$ and $\nabla h = \nabla h^*$ is the identity. Then Theorems 15.1.1 and 15.1.2 reduce to the statement that

$$\mathbb{E}[\|Y - f(X)\|_2^2] = \mathbb{E}[\|Y - \mathbb{E}[Y | X]\|_2^2] + \mathbb{E}[\|\mathbb{E}[Y | X] - f(X)\|_2^2],$$

so we may also see the decompositions of the theorems as bias/variance expansions. \diamond

15.2 Measuring calibration

The first step to building a practicable theory of calibration is to define and measure the calibration of a predictor f . The first step, defining a calibrated predictor, is relatively easy, but measuring how “close” a particular predictor f is to being calibrated raises several challenges, as typical and naive measures of calibration are impossible to estimate. Thus, in this section, we develop several quantities to measure calibration, providing a main theorem relating the different quantities to one another and demonstrating a simple technique to estimate one of them, returning in Section 15.5.2 to show the equivalences between the measures.

We begin with a natural candidate for calibration: the expected difference, or expected calibration error,

$$\text{ece}(f) := \mathbb{E}[|\mathbb{E}[Y | f(X)] - f(X)|]. \quad (15.2.1)$$

The calibration error (15.2.1) is 0 if and only if f is perfectly calibrated, as then $\mathbb{E}[Y | f(X)] = f(X)$, and it is positive otherwise. Unfortunately, while the next lemma guarantees that ece is lower semi-continuous, it is not continuous.

Lemma 15.2.1. *The expected calibration error ece is lower semi-continuous with respect to $L^1(P)$ on \mathcal{F} , that is, if $\mathbb{E}[|f_n(X) - f(X)|] \rightarrow 0$ and $f \in L^1(P)$, then*

$$\liminf_n \text{ece}(f_n) \geq \text{ece}(f).$$

This result requires some delicate measure-theoretic arguments, so we defer it to the technical proofs (see Section 15.6.1). The discontinuity of ece is relatively easy to show, however, even in very simple cases.

Example 15.2.2 (Discontinuity of the calibration error): Let $Y \in \{0, 1\}$ be a Bernoulli random variable, and let $X \in \{0, 1\}$. Take $Y = X$ with probability 1. Then the predictor that always predicts $\frac{1}{2}$ is perfectly calibrated, but if for $\epsilon \in [0, \frac{1}{2}]$ we define f_ϵ by

$$f_\epsilon(0) = \frac{1}{2} - \epsilon \quad \text{and} \quad f_\epsilon(1) = \frac{1}{2} + \epsilon$$

then we see that $\text{ece}(f_\epsilon) = \frac{1}{2} - \epsilon$, while $\text{ece}(f_0) = 0$. Certainly $f_\epsilon \rightarrow f_0$ in any L^p distance on functions, while $\lim_{\epsilon \rightarrow 0} \text{ece}(f_\epsilon) = \frac{1}{2}$. \diamond

15.2.1 The impossibility of measuring calibration

The discontinuity Example 15.2.2 highlights suggests that estimating calibration $\text{ece}(f)$ for a fixed function f should be nontrivial, and indeed, using the tools on functional estimation and testing we develop in Chapter 13, we can show strong lower bounds for estimating the calibration error unless one makes unjustifiable assumptions about the distribution of $Y | f(X)$. The precise reasons differ a bit from the discontinuity of $\text{ece}(f)$ in f , though the intuition is relatively straightforward: if $f(X)$ has a density, then even given a very large sample (X_1^n, Y_1^n) , all the observations $f(X_i)$ will be distinct, and we have no *a priori* reason to assume that $\mathbb{E}[Y | f(X)]$ should be continuous in the predicted value $f(X)$.

To make this more precise, fix a function f whose calibration error we wish to evaluate, and consider a hypothesis test of $H_0 : \text{ece}(f) = 0$ against alternatives that f is miscalibrated, $H_1 : \text{ece}(f) \geq \gamma$ for some $\gamma > 0$. We observe predictions $f(X_i)$ and outcomes Y_i , that is, pairs

$$Z_i = (f(X_i), Y_i)$$

drawn i.i.d.; the coming lower bound holds if $\mathcal{X} = [0, 1]$ and $f(x) = x$, so in many cases, observing X is of no help. Recall the (worst-case) test risk from Section 13.3, that for the testing problem between classes $H_0 : P \in \mathcal{P}_0$ and $H_1 : P \in \mathcal{P}_1$ of distributions,

$$R_n(\Psi \mid \mathcal{P}_0, \mathcal{P}_1) := \sup_{P \in \mathcal{P}_0} P(\Psi(Z_1^n) \neq 0) + \sup_{P \in \mathcal{P}_1} P(\Psi(Z_1^n) \neq 1).$$

Because we consider the function f fixed and ask only whether we can evaluate its calibration error under an (unknown) distribution P , we denote the expected calibration error of f under P via $\text{ece}_P(f) = \mathbb{E}_P[|\mathbb{E}_P[Y \mid f(X)] - f(X)|]$. We thus consider testing perfect calibration $H_0 : \text{ece}(f) = 0$ against alternatives $H_1 : \text{ece}(f) \geq \gamma$ of miscalibration for $\gamma > 0$, defining

$$\mathcal{P}_\gamma = \{\text{distributions } P \text{ on } (X, Y) \mid \text{ece}_P(f) \geq \gamma\}$$

as the collection of distributions for which f is $(\frac{1}{2} - \gamma)$ mis-calibrated.

Theorem 15.2.3. *Let $f : \mathcal{X} \rightarrow [0, 1]$ be a predictor of $Y \in \{0, 1\}$. Assume for some $0 < c < \frac{1}{2}$ that $f(\mathcal{X}) \cap [c, 1 - c]$ has cardinality at least N . Then there is a distribution P_0 such that $\text{ece}_{P_0}(f) = 0$ and for any $0 < \gamma \leq c$,*

$$\inf_{\Psi} R_n(\Psi \mid \{P_0\}, \mathcal{P}_\gamma) \geq 1 - \frac{n\gamma^2}{2\sqrt{N}} \frac{1}{c(1-c)}.$$

Before proving Theorem 15.2.3, we note the following immediate corollary; part (ii) follows from part (i), which follows by taking $N \uparrow \infty$ in the theorem.

Corollary 15.2.4. *Let the conditions of Theorem 15.2.3 hold and let $\mathcal{P}_0 = \{P \mid \text{ece}_P(f) = 0\}$.*

(i) *If there exists $0 < c < \frac{1}{2}$ such that $f(\mathcal{X}) \cap [c, 1 - c]$ has infinite cardinality, then \mathcal{P}_0 is non-empty and for any $0 < \gamma \leq c$,*

$$\liminf_n \inf_{\Psi} R(\Psi \mid \mathcal{P}_0, \mathcal{P}_\gamma) = 1.$$

(ii) *If there exists a neighborhood U of $\frac{1}{2}$ such that $U \subset f(\mathcal{X})$, then \mathcal{P}_0 is non-empty and for any $\gamma < \frac{1}{2}$, the minimax test risk satisfies*

$$\liminf_n \inf_{\Psi} R(\Psi \mid \mathcal{P}_0, \mathcal{P}_\gamma) = 1.$$

In brief, no test exists that is better than random guessing for testing between

$$H_0 : \text{ece}(f) = 0 \quad \text{and} \quad H_1 : \text{ece}(f) \geq c$$

given access to the predictions $f(X_i)$ and observed outcomes Y_i . The theorem and corollary apply to binary prediction models with $Y \in \{0, 1\}$, but the results immediately extend to more complicated prediction problems where Y is vector-valued or multiclass.

Proof The proof relies on the convex hull testing lower bound from Proposition 13.3.1. Without loss of generality, we can assume that $\mathcal{X} \subset [0, 1]$ and that $f(x) = x$ by transforming the input space. Let $S = f(X)$ be the (random) scores that f outputs.

We first construct the perfectly calibrated distribution P_0 and miscalibrated family \mathcal{P}_γ . Define the distribution P_0 so that S is uniform on distinct points $s_1, \dots, s_N \in [c, 1 - c]$ and $Y \mid S = s \sim \text{Bernoulli}(s)$, that is, given $S = s$, $Y = 1$ with probability s and $Y = 0$ with probability $1 - s$. By

construction, $\text{ece}_{P_0}(f) = 0$. To construct the particular members of the alternative family \mathcal{P}_γ , for each $j \in [N]$, define the “tilting” function

$$\phi_j(y, s) := \left(\frac{y}{s_j} - \frac{1-y}{1-s_j} \right) \mathbf{1}\{s = s_j\}.$$

Then $\mathbb{E}_0[\phi_j(Y, S)] = 0$ while

$$\text{Var}_0(\phi_j(Y, S)) = \frac{1}{N} \mathbb{E}_0 \left[\left(\frac{Y}{s_j} - \frac{1-Y}{s_j} \right)^2 \mid S = s_j \right] = \frac{1}{N} \left(\frac{1}{s_j} + \frac{1}{1-s_j} \right) = \frac{1}{N} \frac{1}{s_j(1-s_j)}.$$

Note that $|\phi_j(y, s)| \leq \frac{1}{c}$ as $c < \frac{1}{2}$, and if we define the vector $\phi(y, s) = (\phi_1(y, s), \dots, \phi_N(y, s))$, then $\|\phi(y, s)\|_0 \leq 1$ (that is, the number of non-zero entries is at most 1). Now as $\gamma \in [0, c]$, for each $v \in \{-1, 1\}^N$ we may define the tilted distribution P_v with

$$P_v(Y = y, S = s) = (1 + \gamma \langle v, \phi(y, s) \rangle) P_0(Y = y, S = s),$$

which is a valid distribution whenever $\gamma \leq c$, as $|\langle v, \phi(y, s) \rangle| \leq \frac{1}{c}$. We compute the calibration error for distributions $P \in \{P_v\}$. Noting that S is still uniform on $\{s_1, \dots, s_N\}$ under P_v , we have

$$\mathbb{E}_v[Y \mid S = s_j] = s_j + \gamma v_j \mathbb{E}[\phi_j(Y, s_j) Y \mid S = s_j] = s_j + \gamma v_j,$$

and so $\text{ece}_{P_v}(f) = \frac{1}{N} \sum_{j=1}^N \gamma |v_j| = \gamma$. In particular, we have $P_v \in \mathcal{P}_\gamma$.

Lastly, we compute a bound on the testing error. For this, we recall Lemma 13.2.3. Letting $\overline{P^n} = \frac{1}{2^N} \sum_v P_v^n$, we have

$$\begin{aligned} D_{\chi^2}(\overline{P^n} \| P_0^n) + 1 &= \frac{1}{2^{2N}} \sum_{v, v'} \mathbb{E}_0 [(1 + \gamma \langle v, \phi(Y, S) \rangle)(1 + \gamma \langle v', \phi(Y, S) \rangle)]^n \\ &= \frac{1}{2^{2N}} \sum_{v, v'} \left(1 + \gamma^2 v^\top \text{Cov}_0(\phi(Y, S)) v' \right)^n \end{aligned}$$

because the sampling is i.i.d. By our variance calculation for ϕ and that each ϕ_j has disjoint support, we have $\text{Cov}_0(\phi(Y, S)) = \frac{1}{N} \text{diag}([\frac{1}{s_j(1-s_j)}]_{j=1}^N)$, and so

$$D_{\chi^2}(\overline{P^n} \| P_0^n) + 1 = \mathbb{E} \left[\left(1 + \frac{\gamma^2}{N} \sum_{j=1}^N \frac{V_j V'_j}{s_j(1-s_j)} \right)^n \right] \leq \mathbb{E} \left[\exp \left(\frac{n\gamma^2}{N} \sum_{j=1}^N \frac{V_j V'_j}{s_j(1-s_j)} \right) \right]$$

where the expectation is over $V, V' \stackrel{\text{iid}}{\sim} \text{Uniform}(\{\pm 1\}^N)$. But of course $V_j V'_j$ i.i.d. random signs, and hence 1-sub-Gaussian, so that

$$D_{\chi^2}(\overline{P^n} \| P_0^n) + 1 \leq \exp \left(\frac{n^2 \gamma^4}{N^2} \sum_{j=1}^N \frac{1}{s_j^2(1-s_j)^2} \right) \leq \exp \left(\frac{n^2 \gamma^4}{2N} \frac{1}{c^2(1-c)^2} \right)$$

because $c \leq s_j \leq 1-c$. Apply Proposition 13.3.1 and Pinsker's inequality (Propositions 2.2.8 and 2.2.9) to see that

$$\inf_{\Psi} R(\Psi \mid \{P_0\}, \mathcal{P}_\gamma) \geq 1 - \sqrt{\frac{1}{2} \log(1 + D_{\chi^2}(\overline{P^n} \| P_0))} \geq 1 - \sqrt{\frac{n^2 \gamma^4}{4N} \frac{1}{c^2(1-c)^2}}.$$

Taking square roots gives the result. \square

15.2.2 Alternative calibration measures

The fundamental impossibility results in Theorem 15.2.3 and Corollary 15.2.4, even in the binary prediction case, suggest that we should choose some more easily estimable measure for calibration. In Section 15.5 we provide formal definitions for calibration measures to be continuous (or Lipschitz continuous) and equivalent to one another. Here, we provide the alternative definitions of calibration we consider, giving a corollary that captures their relationships for multiclass classification, and then describing how to estimate one of them. Let us take the general setting of this chapter, where the label space $\mathcal{Y} \subset \mathbb{R}^k$ and P is a distribution on $\mathcal{X} \times \mathcal{Y}$. Let \mathcal{F} be a collection of functions mapping $\mathcal{X} \rightarrow \mathbb{R}^k$ and integrable with respect to P , that is, $\mathbb{E}[\|f(X)\|] < \infty$ for each $f \in \mathcal{F}$.

In brief, we require that a calibration measure $M : \mathcal{F} \rightarrow \mathbb{R}_+$ be *sound* (in analogy with proof systems, where soundness means nothing false can be proved), meaning that

$$M(f) = 0 \text{ implies } \mathbb{E}[Y | f(X)] = f(X) \quad (15.2.2a)$$

and *complete* (continuing the analogy, that everything true can be proved), meaning that

$$\mathbb{E}[Y | f(X)] = f(X) \text{ implies } M(f) = 0. \quad (15.2.2b)$$

We begin by considering types of *distance to calibration*. Let $\mathcal{C}(P)$ denote those functions g that are perfectly calibrated for P , that is, $\mathcal{C}(P) = \{g : \mathcal{X} \rightarrow \mathbb{R}^k \mid \mathbb{E}_P[Y | g(X)] = g(X)\}$ (where the defining equality holds with P -probability 1 over X). The set \mathcal{P} always consists at least of the constant function $g(X) = \mathbb{E}_P[Y]$ and so is non-empty (but is typically larger). Then we call the minimum $L^1(P)$ distance of a function f to the set $\mathcal{C}(P)$ the *distance to calibration*

$$d_{\text{cal}}(f) := \inf_g \{\mathbb{E}[\|g(X) - f(X)\|] \text{ s.t. } g \in \mathcal{C}(P)\}. \quad (15.2.3)$$

It is not always clear how to estimate the distance $d_{\text{cal}}(f)$, making using it sometimes challenging.

We also consider a complementary quantity that relies on an alternative variational characterization. Let $\mathcal{W} \subset \{\mathbb{R}^k \rightarrow \mathbb{R}^k\}$ be a symmetric collection of functions, meaning that $w \in \mathcal{W}$ implies $-w \in \mathcal{W}$. We can view any such collection as potential witnesses of miscalibration, in that

$$\mathbb{E}[\langle w(f(X)), Y - f(X) \rangle] = \mathbb{E}[\langle w(f(X)), \mathbb{E}[Y | f(X)] - f(X) \rangle]$$

and so if w can “witness” the portions of space where $f(X) \not\approx \mathbb{E}[Y | f(X)]$, it can certify miscalibration. We then arrive at what we term the *calibration error relative to the class \mathcal{W}* ,

$$\text{CE}(f, \mathcal{W}) := \sup_{w \in \mathcal{W}} \mathbb{E}[\langle w(f(X)), Y - f(X) \rangle]. \quad (15.2.4)$$

Depending on the class \mathcal{W} , this is sometimes called the *weak calibration error*, and with large enough classes, we can recover the classical expected calibration error (15.2.1).

Example 15.2.5 (Recovering expected calibration error): For a norm $\|\cdot\|$, let the set \mathcal{W} be the collection of all functions w with bound $\sup_s \|w(s)\|_* \leq 1$. Then

$$\text{CE}(f, \mathcal{W}) = \mathbb{E} \left[\sup_{\|w\|_* \leq 1} \langle w, \mathbb{E}[Y | f(X)] - f(X) \rangle \right] = \mathbb{E}[\|\mathbb{E}[Y | f(X)] - f(X)\|] = \text{ece}(f),$$

the expected calibration error. \diamond

It is more interesting to consider restricted classes; one of particular interest to us is that of bounded Lipschitz functions. Let

$$\mathcal{W}_{\|\cdot\|} := \left\{ w : \mathbb{R}^k \rightarrow \mathbb{R}^k \mid \|w(s_0) - w(s_1)\|_* \leq \|s_0 - s_1\| \text{ and } \|w(s)\|_* \leq 1 \text{ for all } s, s_0, s_1 \right\} \quad (15.2.5)$$

denote the collection of functions bounded by 1 in $\|\cdot\|_*$ and that are 1-Lipschitz with respect to $\|\cdot\|$. Then (as we see presently) we can at least estimate the calibration error relative to the class \mathcal{W} in the definition (15.2.4).

The final calibration measure we consider repos on the idea of quantizing or partitioning the output space, which relates to the idea of “binning” predictions that the literature on calibration frequently considers. Here, we consider averages of Y conditioned on predictions in larger sets. Thus, instead of evaluating the precise conditioning $\mathbb{E}[Y \mid f(X)]$ we to look instead at the expectation of Y conditional on $f(X) \in A$ for a set A , so that a predicted score is (nearly) calibrated if the diameter $\text{diam}(A)$ is small, and $\mathbb{E}[Y \mid f(X) \in A] \approx s$ for some $s \in A$. Given a partition \mathcal{A} of the space $\mathcal{M} = \text{Conv}(\mathcal{Y})$, it is then natural to evaluate the average error for each element of A (weighting by the probability of A), and consider the calibration error (15.2.4) for indicator functions of $A \in \mathcal{A}$, where we abuse notation slightly to define

$$\text{CE}(f, \mathcal{A}) := \sum_{A \in \mathcal{A}} \|\mathbb{E}[(f(X) - Y)\mathbf{1}\{f(X) \in A\}]\| = \sum_{A \in \mathcal{A}} \|\mathbb{E}[f(X) - Y \mid f(X) \in A]\| \mathbb{P}(f(X) \in A).$$

Indeed, taking a supremum over all such partitions gives $\sup_{\mathcal{A}} \text{CE}(f \mid \mathcal{A}) = \mathbb{E}[\|\mathbb{E}[Y \mid f(X)] - f(X)\|]$, the original expected calibration error (15.2.1). Additionally, and here we elide details, if $f(X)$ is a continuous random variable with suitably nice density and \mathcal{A}_n denotes any partition satisfying $\text{diam}(A) \leq 1/n$ for $A \in \mathcal{A}_n$, then $\lim_n \text{CE}(f, \mathcal{A}_n) = \mathbb{E}[\|\mathbb{E}[Y \mid f(X)] - f(X)\|]$. Instead of considering $\text{CE}(f, \mathcal{A})$ directly, we optimize over all partitions, but penalize the average size of elements of \mathcal{A} , giving the *partitioned calibration error*

$$\text{pce}(f) := \inf_{\mathcal{A}} \left\{ \text{CE}(f, \mathcal{A}) + \sum_{A \in \mathcal{A}} \text{diam}(A) \mathbb{P}(f(X) \in A) \right\}. \quad (15.2.6)$$

Each of these is equivalent to within polynomial scaling.

Corollary 15.2.6. *Let $\mathcal{Y} \subset \mathbb{R}^k$ have finite diameter and $\|\cdot\|$ be any norm. Then each of the calibration measures d_{cal} , $\text{CE}(\cdot, \mathcal{W}_{\|\cdot\|})$, and pce in definitions (15.2.3), (15.2.4), and (15.2.6) is sound and complete (15.2.2). Additionally, let $\mathcal{Y} = \{e_1, \dots, e_k\}$ and $\|\cdot\| = \|\cdot\|_1$ be the ℓ_1 -norm. Then for any $f : \mathcal{X} \rightarrow \mathcal{M} = \text{Conv}(\mathcal{Y})$,*

$$\frac{1}{2} \text{CE}(f, \mathcal{W}_{\|\cdot\|}) \leq d_{\text{cal}}(f) \leq \text{CE}(f, \mathcal{W}_{\|\cdot\|}) + 2\sqrt{k \text{CE}(f, \mathcal{W}_{\|\cdot\|})}$$

and

$$d_{\text{cal}}(f) \leq \text{pce}(f) \leq d_{\text{cal}}(f) + 2\sqrt{k d_{\text{cal}}(f)}.$$

Corollary 15.2.6 will come as a consequence of the deeper development we pursue in Section 15.5.

Here, we take Corollary 15.2.6 as motivation to give the type of typical result that justifies calibration estimates. As any of the calibration measures is roughly equivalent (except ece), measuring any of them on a sample can provide evidence for or against calibration of a predictor f . We focus

on the simpler binary case in which $f : \mathcal{X} \rightarrow [0, 1]$ and let \mathcal{W}_{Lip} be bounded Lipschitz functions $w : [0, 1] \rightarrow [-1, 1]$. Given a sample (X_1^n, Y_1^n) , the empirical variant of $\text{CE}(f, \mathcal{W})$ is

$$\widehat{\text{CE}}_n(f) := \sup_{\|w\|_\infty \leq 1} \left\{ \frac{1}{n} \sum_{i=1}^n w_i(Y_i - f(X_i)) \text{ s.t. } |w_i - w_j| \leq |f(X_i) - f(X_j)| \text{ for } i, j \leq n \right\}.$$

By combining uniform covering bounds for the class of Lipschitz functions with a standard concentration inequality, we then have the following convergence guarantee for $\widehat{\text{CE}}_n$.

Proposition 15.2.7. *There exists a numerical constant C such that for any $\delta > 0$,*

$$\left| \widehat{\text{CE}}_n(f) - \text{CE}(f, \mathcal{W}_{\text{Lip}}) \right| \leq C \frac{\sqrt{\log \frac{n}{\delta}}}{n^{1/3}}$$

with probability at least $1 - \delta$.

Proof Fix $\epsilon > 0$ and let $\mathcal{N}(\epsilon)$ be a minimal ϵ -cover of the set \mathcal{W}_{Lip} in uniform norm, meaning that $\|w - w^{(j)}\|_\infty \leq \epsilon$ for each $w^{(j)} \in \mathcal{N}(\epsilon)$, and let $N(\epsilon)$ be its (minimal) cardinality. Then $\log N(\epsilon) \lesssim \frac{1}{\epsilon} \log \frac{1}{\epsilon}$ (recall Proposition 10.2.3 and Eq. (10.2.4)). For shorthand, let the error vector $E \in [-1, 1]^n$ have entries $E_i = Y_i - f(X_i)$, and abusing notation, for $w \in \mathcal{W}_{\text{Lip}}$ let $\langle w, E \rangle_n = \frac{1}{n} \sum_{i=1}^n w(f(X_i)) E_i$. Then for any $w \in \mathcal{W}_{\text{Lip}}$, there exists $i \leq N(\epsilon)$ such that

$$|\langle w, E \rangle_n - \langle w^{(i)}, E \rangle_n| \leq \epsilon,$$

while $\widehat{\text{CE}}_n(f) = \sup_{w \in \mathcal{W}_{\text{Lip}}} \langle w, E \rangle_n$. In particular, we have

$$\left| \widehat{\text{CE}}_n(f) - \text{CE}(f, \mathcal{W}_{\text{Lip}}) \right| \leq \sup_{w \in \mathcal{W}_{\text{Lip}}} |\langle w, E \rangle_n - \mathbb{E}[\langle w, E \rangle_n]| \leq \max_{w \in \mathcal{N}(\epsilon)} |\langle w, E \rangle_n - \mathbb{E}[\langle w, E \rangle_n]| + 2\epsilon.$$

Thus for any $t \geq 0$, we have

$$\begin{aligned} \mathbb{P} \left(\left| \widehat{\text{CE}}_n(f) - \text{CE}(f, \mathcal{W}_{\text{Lip}}) \right| \geq t \right) &\leq \mathbb{P} \left(\max_{w \in \mathcal{N}(\epsilon)} |\langle w, E \rangle_n - \mathbb{E}[\langle w, E \rangle_n]| \geq t - 2\epsilon \right) \\ &\leq 2N(\epsilon) \exp \left(-\frac{n[t - 2\epsilon]_+^2}{2} \right) \end{aligned}$$

by the Azuma-Hoeffding inequality and a union bound. Take $\epsilon = n^{-1/3}$ and $t = Cn^{-1/3} \sqrt{\log \frac{n}{\delta}}$ for an appropriate numerical constant C to obtain the proposition. \square

Summarizing, while the expected calibration error is fundamentally inestimable, there are alternative measures that are both sound and complete, and they can admit reasonable estimators. As the class size k grows, however, it can become statistically infeasible to estimate the calibration of predictors f , so that one must consider alternative metrics. The exercises and bibliography explore these questions in more detail.

15.3 Auditing and improving calibration at the population level

Theorems 15.1.1 and 15.1.2 provide decompositions of the expected loss of a predictor

$$\mathbb{E}[\ell(f(X), Y)] = \mathbb{E}[\ell(\mathbb{E}[Y | f(X)], Y)] + \mathbb{E}[D_h(\mathbb{E}[Y | f(X)], f(X))]$$

into an average loss and an expected divergence between $f(X)$ and $\mathbb{E}[Y | f(X)]$, where h is the negative (generalized) entropy (14.1.6) associated with the loss ℓ , so that the loss has representation $\ell(\mu, y) = -h(\mu) - \langle \nabla h(\mu), y - \mu \rangle$. This suggests an approach to improving a predictor $f : \mathcal{X} \rightarrow \mathbb{R}^k$ without compromising its average loss: make it closer to being calibrated, so that $\mathbb{E}[Y | f(X)] \approx f(X)$. Here, we make this idea precise by using the weak calibration (15.2.4): if there exists a witness function w certifying that $\mathbb{E}[\langle w(f(X)), Y - f(X) \rangle] \gg 0$, then we can post-process f to $f(X) + \eta w(f(X))$ for some stepsize $\eta > 0$ and *only* improve the expected loss. We first develop the idea in the context of the squared error, where the calculations are cleanest, and extend it to general proper losses based on convex conjugates (as in Section 14.3) immediately after. Combining the ideas we develop, we also provide a (population-level) algorithm to transform a function f by post-processing its outputs that guarantees the result is nearly calibrated relative to a class \mathcal{W} of witnesses. This provides an algorithmic proof quantitatively relating the calibration error $\text{CE}(f, \mathcal{W})$ relative to a class \mathcal{W} to the improvement achievable in minimizing $\mathbb{E}[\ell(f(X), Y)]$ by post-composition $g \circ f$.

15.3.1 The post-processing gap and calibration audits for squared error

Consider a thought experiment: instead of using f to make predictions, we use a postprocessing $g \circ f$, where $g : \mathbb{R}^k \rightarrow \mathbb{R}^k$ has the (suggestively chosen) form $g(v) = v + w(v)$, where $w(v) = (g(v) - v)$. Then using the representation $\ell(\mu, y) = -h(\mu) - \langle \nabla h(\mu), y - \mu \rangle$ for the proper loss, we recall Theorem 15.1.1 and for $\mu(f(X)) := \mathbb{E}[Y | f(X)]$ expand

$$\begin{aligned} \mathbb{E}[\ell(g \circ f(X), Y)] &= \mathbb{E}[-h(g \circ f(X)) - \langle \nabla h(g \circ f(X)), Y - g \circ f(X) \rangle] \\ &= \mathbb{E}[-h(\mu(f(X)))] + \mathbb{E}[h(\mu(f(X))) - \langle \nabla h(g \circ f(X)), Y - g \circ f(X) \rangle] \\ &= \mathbb{E}[\ell(\mathbb{E}[Y | f(X)], Y)] + \mathbb{E}[D_h(\mathbb{E}[Y | f(X)], g \circ f(X))], \end{aligned}$$

where the final equality uses the linearity of $y \mapsto \ell(\mu, y)$, that is,

$$\mathbb{E}[\ell(g \circ f(X), Y)] = \mathbb{E}[\ell(\mathbb{E}[Y | f(X)], Y)] + \mathbb{E}[D_h(\mathbb{E}[Y | f(X)], f(X) + w(f(X)))]. \quad (15.3.1)$$

We have decomposed the expected loss $\mathbb{E}[\ell(g \circ f(X), Y)]$ into a term that post-processing does not change, which measures the sharpness with which $\mathbb{E}[Y | f(X)]$ predicts Y , and a divergence term D_h measuring the error in calibration of $g \circ f(X) = f(X) + w(f(X))$ for $\mathbb{E}[Y | f(X)]$.

The expansion (15.3.1) points toward an ability to postprocess *any* prediction function $f : \mathcal{X} \rightarrow \mathbb{R}^k$ to both (i) obtain calibration relative to a class of functions \mathcal{W} , as in Definition (15.2.4), and (ii) improve the expected loss $\mathbb{E}[\ell(f(X), Y)]$. Moreover, this improvement is monotone, in that changes “toward” calibration guarantee smaller expected loss, an improvement over the less refined results in Theorems 15.1.1 and 15.1.2. To that end, define the *post-processing gap* for the (proper) loss ℓ and function f relative to the class \mathcal{W} of functions $\mathbb{R}^k \rightarrow \mathbb{R}^k$ by

$$\text{gap}(\ell, f, \mathcal{W}) := \mathbb{E}[\ell(f(X), Y)] - \inf_{w \in \mathcal{W}} \mathbb{E}[\ell(f(X) + w(f(X)), Y)]. \quad (15.3.2)$$

The gap (15.3.2) is fundamentally tied to the calibration error relative to the class \mathcal{W} .

We specialize here to the simpler case of the squared error, as the statements are most transparent. We focus exclusively on symmetric convex collections of functions \mathcal{W} , meaning that if $w \in \mathcal{W}$, then $-w \in \mathcal{W}$, and \mathcal{W} is convex.

Proposition 15.3.1. *Let $\ell(\mu, y) = \frac{1}{2} \|y - \mu\|_2^2$ be the squared error (Brier score), and let \mathcal{W} be a symmetric convex collection of functions, each 1-Lipschitz with respect to the ℓ_2 -norm $\|\cdot\|_2$. Define $R^2(f) = \sup_{w \in \mathcal{W}} \mathbb{E}[\|w(f(X))\|_2^2]$. Then*

$$\frac{1}{2} \min \left\{ \text{CE}(f, \mathcal{W}), \frac{\text{CE}(f, \mathcal{W})^2}{R^2(f)} \right\} \leq \text{gap}(\ell, f, \mathcal{W}) \leq \text{CE}(f, \mathcal{W})$$

Proof Fix x and let $\mu = \mathbb{E}[Y \mid f(X) = f(x)] \in \text{Conv}(\mathcal{Y})$ and $w = w(f(x))$ be a potential update to $f(x)$. Then because $\ell(\mu, y) = \frac{1}{2} \|\mu - y\|_2^2$, for any $y \in \mathcal{Y}$

$$\ell(\mu, y) + \langle \nabla \ell(\mu, y), w \rangle + \frac{1}{2} \|w\|^2 = \ell(\mu + w, y).$$

Recognizing that $\nabla \ell(\mu, y) = (\mu - y)$, for any $w \in \mathcal{W}$ we therefore have

$$\begin{aligned} -\mathbb{E}[\langle f(X) - Y, w(f(X)) \rangle] - \frac{1}{2} \mathbb{E}[\|w(f(X))\|_2^2] &\leq \mathbb{E}[\ell(f(X), Y)] - \mathbb{E}[\ell(f(X) + w(f(X)), Y)] \\ &\leq -\mathbb{E}[\langle f(X) - Y, w(f(X)) \rangle]. \end{aligned}$$

Taking suprema over w on each side of the preceding inequalities and using the symmetry of \mathcal{W} gives

$$\begin{aligned} \sup_{w \in \mathcal{W}} \left\{ \mathbb{E}[\langle f(X) - Y, w(f(X)) \rangle] - \frac{1}{2} \mathbb{E}[\|w(f(X))\|_2^2] \right\} &\leq \text{gap}(\ell, f, \mathcal{W}) \\ &\leq \sup_{w \in \mathcal{W}} \mathbb{E}[\langle f(X) - Y, w(f(X)) \rangle]. \end{aligned}$$

Because $\text{CE}(f, \mathcal{W}) = \sup_{w \in \mathcal{W}} \mathbb{E}[\langle f(X) - Y, w(f(X)) \rangle]$, we can use the convexity of \mathcal{W} and the definition $R^2(f) := \sup_{w \in \mathcal{W}} \mathbb{E}[\|w(f(X))\|_2^2]$ to see that for any $\eta \in [0, 1]$, we may replace w with $\eta \cdot w \in \mathcal{W}$, and we have

$$\sup_{\eta \in [0, 1]} \left[\eta \text{CE}(f, \mathcal{W}) - \frac{\eta^2}{2} R^2(f) \right] \leq \text{gap}(\ell, f, \mathcal{W}) \leq \text{CE}(f, \mathcal{W}).$$

Maximizing over η on the left side, we choose $\eta = \min\{1, \frac{\text{CE}(f, \mathcal{W})}{R^2(f)}\}$ to obtain the proposition. \square

As an immediate corollary, we see that if $\mathcal{W} = \mathcal{W}_{\|\cdot\|_2}$ consists of the 1-Lipschitz functions with $\|w(\cdot)\|_2 \leq 1$, we have a cleaner guarantee.

Corollary 15.3.2. *Let $\mathcal{W} = \mathcal{W}_{\|\cdot\|_2}$ and the conditions of Proposition 15.3.1 hold. Then*

$$\frac{1}{2 \text{diam}(\mathcal{Y})^2} \text{CE}(f, \mathcal{W})^2 \leq \text{gap}(\ell, f, \mathcal{W}) \leq \text{CE}(f, \mathcal{W}).$$

Thus, the calibration error upper and lower bounds the gap between the expected loss of f and a post-processed version of f . This yields a nearly operational interpretation of the calibration error relative to the class \mathcal{W} : it is, to within a square, exactly the amount we could improve the expected loss of the function f by postprocessing f itself.

15.3.2 Calibration audits for losses based on conjugate linkages

Recall as in Section 14.3.1 that, by a transformation tied to the loss ℓ via its associated generalized negative entropy, we may define the surrogate

$$\varphi(s, y) := h^*(s) - \langle s, y \rangle,$$

and we may transform arbitrary scores $s \in \mathbb{R}^k$ to predictions via the conjugate link (14.3.1), that is,

$$\text{pred}_h(s) = \underset{\mu}{\operatorname{argmin}} \{-\langle s, \mu \rangle + h(\mu)\} = \nabla h^*(s).$$

So long as h is appropriately smooth, these satisfy $\ell(\text{pred}_h(s), y) = \varphi(s, y)$. In complete analogy with the post-processing gap (15.3.2) when we assume f makes predictions in (the affine hull of) \mathcal{Y} , we can define the *surrogate post-processing gap*

$$\text{gap}(\varphi, f, \mathcal{W}) := \mathbb{E}[\varphi(f(X), Y)] - \inf_{w \in \mathcal{W}} \mathbb{E}[\varphi(f(X) + w(f(X)), Y)]. \quad (15.3.3)$$

In spite of the similarity with definition (15.3.2), the actual predictions of Y from f in this case come via the link $\text{pred}_h(f(X))$. Thus, in this case we instead consider the calibration error relative to a class \mathcal{W} but after the composition of f with $\text{pred}_h = \nabla h^*$, so that

$$\text{CE}(\text{pred}_h \circ f, \mathcal{W}) = \sup_{w \in \mathcal{W}} \mathbb{E}[\langle w(f(X)), Y - \text{pred}_h(f(X)) \rangle] = \sup_{w \in \mathcal{W}} \mathbb{E}[\langle w(f(X)), Y - \nabla h^*(f(X)) \rangle],$$

where as always we assume that the class of witness functions satisfies $\mathcal{W} = -\mathcal{W}$. When the prediction function is continuous enough in s , we can give an analogue of Proposition 15.3.1 to the more general surrogate case. To that end, we assume that the conjugate h^* has Lipschitz continuous gradient with respect to the dual norm $\|\cdot\|_*$, meaning that

$$\|\nabla h^*(s_0) - \nabla h^*(s_1)\| \leq \|s_0 - s_1\|_*$$

for all $s_0, s_1 \in \mathbb{R}^k$. This is equivalent (see Proposition C.2.6) to the negative entropy h being strongly convex with respect to the norm $\|\cdot\|$, and also immediately implies that

$$\varphi(s + w, y) \leq \varphi(s, y) + \langle \nabla_s \varphi(s, y), w \rangle + \frac{\|w\|_*^2}{2}.$$

Example 15.3.3 (Multiclass logistic regression): For multiclass logistic regression, where we take $h(p) = \sum_{j=1}^k p_j \log p_j$, we know that h is strongly convex with respect to the ℓ_1 norm (this is Pinsker's inequality; see inequality (2.2.11)). Thus, the conjugate $h^*(s) = \log(\sum_{j=1}^k e^{s_j})$ has Lipschitz gradient with respect to the ℓ_∞ norm, meaning that for the prediction link

$$\text{pred}_h(s) = \left[\frac{e^{s_y}}{\sum_{j=1}^k e^{s_j}} \right]_{y=1}^k,$$

we have

$$\|\text{pred}_h(s) - \text{pred}_h(s')\|_1 \leq \|s - s'\|_\infty$$

for all $s, s' \in \mathbb{R}^k$. \diamond

Example 15.3.4 (The squared error): When we measure the error of a predictions in \mathbb{R}^k by the squared ℓ_2 -norm $\frac{1}{2} \|f(x) - y\|_2^2$, this corresponds to the generalized negative entropy $h(\mu) = \frac{1}{2} \|\mu\|_2^2$. In this case, the norm $\|\cdot\| = \|\cdot\|_2 = \|\cdot\|_*$, and we have the self duality $h^* = h$, so that the prediction mapping pred_h is the identity. \diamond

With these examples as motivation, we then have the following generalization of Proposition 15.3.1.

Proposition 15.3.5. *Let the negative generalized entropy h be strongly convex with respect to the norm $\|\cdot\|$ and consider surrogate loss $\varphi(s, y) = h^*(s) - \langle s, y \rangle$. Define $R_*^2(f) := \sup_{w \in \mathcal{W}} \mathbb{E}[\|w(f(X))\|_*^2]$. Then*

$$\frac{1}{2} \min \left\{ \text{CE}(\text{pred}_h \circ f, \mathcal{W}), \frac{\text{CE}(\text{pred}_h \circ f, \mathcal{W})^2}{R_*^2(f)} \right\} \leq \text{gap}(\varphi, f, \mathcal{W}) \leq \text{CE}(\text{pred}_h \circ f, \mathcal{W}).$$

Proof Fix x and let $s = f(x)$ and $w = w(f(x))$, and notice that for any y we have

$$\varphi(s, y) + \langle \nabla \varphi(s, y), w \rangle \leq \varphi(s + w, y) \leq \varphi(s, y) + \langle \nabla \varphi(s, y), w \rangle + \frac{1}{2} \|w\|_*^2.$$

Recognizing that $\nabla \varphi(s, y) = \nabla h^*(s) - y$, for any $w \in \mathcal{W}$ we have

$$\begin{aligned} -\mathbb{E}[\langle \nabla \varphi(f(X), Y), w(f(X)) \rangle] - \frac{1}{2} \mathbb{E}[\|w(f(X))\|_*^2] &\leq \mathbb{E}[\varphi(f(X), Y)] - \mathbb{E}[\varphi(f(X) + w(f(X)), Y)] \\ &\leq -\mathbb{E}[\langle \nabla \varphi(f(X), Y), w(f(X)) \rangle]. \end{aligned}$$

Taking suprema over w on each side and using the symmetry of \mathcal{W} gives

$$\begin{aligned} \sup_{w \in \mathcal{W}} \left\{ \mathbb{E}[\langle \nabla h^*(f(X)) - Y, w(f(X)) \rangle] - \frac{1}{2} \mathbb{E}[\|w(f(X))\|_*^2] \right\} &\leq \text{gap}(\varphi, f, \mathcal{W}) \\ &\leq \sup_{w \in \mathcal{W}} \mathbb{E}[\langle \nabla h^*(f(X)) - Y, w(f(X)) \rangle]. \end{aligned}$$

Because $\text{CE}(\text{pred}_h \circ f, \mathcal{W}) = \sup_{w \in \mathcal{W}} \mathbb{E}[\langle \nabla h^*(f(X)) - Y, w(f(X)) \rangle]$, we can use the convexity of \mathcal{W} and the definition $R_*^2(f) := \sup_{w \in \mathcal{W}} \mathbb{E}[\|w(f(X))\|_*^2]$, to see that for any $\eta \in [0, 1]$, we may replace w with $\eta \cdot w \in \mathcal{W}$ and

$$\sup_{\eta \in [0, 1]} \left[\eta \text{CE}(\text{pred}_h \circ f, \mathcal{W}) - \frac{\eta^2}{2} R_*^2(f) \right] \leq \text{gap}(\varphi, f, \mathcal{W}) \leq \text{CE}(\text{pred}_h \circ f, \mathcal{W}).$$

Set $\eta = \min\{1, \frac{\text{CE}(\text{pred}_h \circ f, \mathcal{W})}{R_*^2(f)}\}$. \square

A corollary specializing to the case of bounded witness functions allows a somewhat cleaner statement, in analogy with Corollary 15.3.2. It provides the same operational interpretation: the calibration error $\text{CE}(f, \mathcal{W})$ of f relative to \mathcal{W} upper and lower bounds improvement possible through postprocessing f .

Corollary 15.3.6. *Let the conditions of Proposition 15.3.5 hold, and additionally assume that the witness functions \mathcal{W} satisfy $\|w(s)\|_* \leq 1$ for all $s \in \mathbb{R}^k$. Then*

$$\frac{1}{2 \text{diam}(\text{dom } h)^2} \text{CE}(\text{pred}_h \circ f, \mathcal{W})^2 \leq \text{gap}(\varphi, f, \mathcal{W}) \leq \text{CE}(\text{pred}_h \circ f, \mathcal{W}).$$

We can give an alternative perspective for this section by focusing on the definitions (15.3.2) and (15.3.3) of the post-processing gap. Suppose we have a proper loss ℓ and we wish to improve the expected loss of a predictor f by post-processing f . When there is little to be gained by replacing f with an adjusted version $f(x) + w(f(x))$ for some $w \in \mathcal{W}$, then f *must be calibrated* with respect to the class \mathcal{W} . So, for example, for a surrogate φ , the function f (really, its associated prediction function $\text{pred}_h \circ f$) is calibrated with respect to \mathcal{W} if and only if $\mathbb{E}[\varphi(f(X) + w(f(X)), Y)] \leq \mathbb{E}[\varphi(f(X), Y)]$ for all $w \in \mathcal{W}$.

As a particular special case to close this section, the standard multiclass logistic loss provides a clean example.

Example 15.3.7 (Multiclass logistic losses, continued): Let h be the negative entropy $h(p) = \sum_{j=1}^k p_j \log p_j$ restricted to the probability simplex $\Delta_k = \{p \in \mathbb{R}_+^k \mid \langle \mathbf{1}, p \rangle = 1\}$ and the surrogate $\varphi(s, y) = \log(\sum_{j=1}^k e^{s_j}) - s_y$. Then for any class \mathcal{W} consisting of functions with $\|w(s)\|_\infty \leq 1$ for all $s \in \mathbb{R}^k$ and any function $f : \mathcal{X} \rightarrow \mathbb{R}^k$,

$$\frac{1}{2} \text{CE}(\text{pred}_h \circ f, \mathcal{W})^2 \leq \mathbb{E}[\varphi(f(X), Y)] - \inf_{w \in \mathcal{W}} \mathbb{E}[\varphi(f(X) + w(f(X)), Y)].$$

(Note that $\text{dom } h$ has diameter 1 in the ℓ_1 -norm.) \diamond

15.3.3 A population-level algorithm for calibration

Implicit in each of the calibration gap bounds in Propositions 15.3.1 and 15.3.5 is bound on the improvement of a predictor f relative to processing outputs with a class \mathcal{W} of functions. This suggests an algorithm for updating the predictions of f to make them calibrated, after which no improvement is possible. While we work at the population level here, similar procedures can allow calibration given access to additional data.

Working in the more general setting of surrogate losses based on the generalized negative entropy h , as these include the standard squared error as a special case, the key idea is that if we find the witness w maximizing $\mathbb{E}[\langle w(f(X)), Y - \text{pred}_h(f(X)) \rangle]$ we can update f with $f - \eta \cdot w \circ f$ for some stepsize η , thus improving the calibration of f relative to the class \mathcal{W} of potential witnesses. In Figure 15.1, we present a prototypical algorithm for achieving this.

The following theorem bounds the convergence of the algorithm.

Theorem 15.3.8. Assume that the surrogate loss φ is nonnegative and that the class of witnesses \mathcal{W} satisfies $R_* := \sup_s \|w(s)\|_* < \infty$. Then the algorithm in Figure 15.1 guarantees that

$$\min_{\tau \leq t} \text{CE}(\text{pred}_h \circ f_\tau, \mathcal{W}) \leq \frac{\sqrt{2R_*^2 \mathbb{E}[\varphi(f_0(X), Y)]}}{\sqrt{t}},$$

and in particular terminates with $\text{CE}(\text{pred}_h \circ f_t, \mathcal{W}) \leq \epsilon$ for some t with

$$t \leq \frac{2R_*^2 \mathbb{E}[\varphi(f_0(X), Y)]}{\epsilon^2}.$$

Proof We begin by showing a one-step progress guarantee beginning from a fixed function f . For any $w : \mathbb{R}^k \rightarrow \mathbb{R}^k$ and any f , we have

$$\mathbb{E}[\varphi(f(X) + \eta w(f(X)), Y)] \leq \mathbb{E}[\varphi(f(X), Y)] + \eta \mathbb{E}[\langle w(f(X)), \nabla h^*(f(X)) - Y \rangle] + \frac{\eta^2}{2} \mathbb{E}[\|w(f(X))\|_*^2].$$

Input: Population distribution P , collection of bounded witness functions \mathcal{W} , generalized negative entropy h strongly convex w.r.t. norm $\|\cdot\|$, initial predictor $f_0 : \mathcal{X} \rightarrow \mathbb{R}^k$, calibration tolerance $\epsilon > 0$

Initialize: set $R_*^2 := \sup_s \|w(s)\|_*^2$

Repeat: for $t = 0, 1, \dots$

i. Find witness w_t maximizing $\mathbb{E}[\langle w(f_t(X)), Y - \text{pred}_h(f_t(X)) \rangle]$

ii. Set $\eta_t = \frac{\mathbb{E}[\langle w_t(f_t(X)), Y - \text{pred}_h(f_t(X)) \rangle]}{R_*^2}$

iii. Update $f_{t+1} = f_t - \eta_t \cdot w_t \circ f_t$

iv. Terminate if

$$\text{CE}(\text{pred}_h \circ f_t, \mathcal{W}) \leq \epsilon.$$

Figure 15.1: Improving calibration relative to the class \mathcal{W}

Let w maximize $\mathbb{E}[\langle w(f(X)), \nabla h^*(f(X)) - Y \rangle]$, so that

$$\mathbb{E}[\langle \varphi(f(X) - \eta w(f(X))), Y \rangle] \leq \mathbb{E}[\langle \varphi(f(X)), Y \rangle] - \eta \text{CE}(\text{pred}_h \circ f, \mathcal{W}) + \frac{\eta^2}{2} R_*^2.$$

Choose $\eta_f = \frac{\text{CE}(\text{pred}_h \circ f, \mathcal{W})}{R_*^2}$ to obtain

$$\mathbb{E}[\langle \varphi(f(X) - \eta_f w(f(X))), Y \rangle] \leq \mathbb{E}[\langle \varphi(f(X)), Y \rangle] - \frac{1}{2} \frac{\text{CE}(\text{pred}_h \circ f, \mathcal{W})^2}{R_*^2}. \quad (15.3.4)$$

Now we apply the obvious inductive argument. Let f_t be a function in the iteration of Algorithm 15.1. Then inequality (15.3.4) guarantees that if $\delta_t^2 := \frac{1}{2} \frac{\text{CE}(\text{pred}_h \circ f_t, \mathcal{W})^2}{R_*^2}$, then

$$\mathbb{E}[\langle \varphi(f_{t+1}(X)), Y \rangle] \leq \mathbb{E}[\langle \varphi(f_t(X)), Y \rangle] - \delta_t^2.$$

In particular,

$$0 \leq \mathbb{E}[\langle \varphi(f_t(X)), Y \rangle] \leq \mathbb{E}[\langle \varphi(f_0(X)), Y \rangle] - \sum_{\tau=0}^{t-1} \delta_\tau^2.$$

In particular,

$$t \min_{\tau < t} \delta_\tau^2 \leq \mathbb{E}[\langle \varphi(f_0(X)), Y \rangle],$$

so that $\min_{\tau < t} \delta_\tau \leq \sqrt{\mathbb{E}[\langle \varphi(f_0(X)), Y \rangle] / t}$. Replacing δ_τ with its definition gives the theorem. \square

15.4 Calibrating: improving squared error by calibration

Sections 15.1 and 15.3 show that at least at the population level, taking a predictor f and modifying (or postprocessing) it to guarantee its calibration can only improve the losses it suffers, whether

those are squared error or general proper losses. That is, by calibrating we can beat (and hence, calibrate) a given predictor. These arguments have exclusively been at the population level, leaving it unclear whether this approach might actually work given a finite sample. While employing these ideas for general losses and general decision settings, where we only guarantee $\mathcal{Y} \subset \mathbb{R}^k$, is challenging because of dimensionality issues, here we show how to improve calibration in finite samples while simultaneously losing little in squared error for binary predictions with $Y \in \{0, 1\}$. That is, we have *calibrating*: from any potential predictor f , we can construct a predictor g with both small calibration error and with (asymptotically) no larger squared error than f , realizing Theorem 15.1.1 but in finite samples.

Let $f : \mathcal{X} \rightarrow [0, 1]$ be any predictor of $Y \in \{0, 1\}$, and consider the squared error loss $\ell(s, y) = (s - y)^2$ with population loss $L(f) = \mathbb{E}[(Y - f(X))^2]$. The idea to improve calibration of f without losing much in accuracy (squared error) is fairly straightforward: we discretize f by binning its predictions so that the number of X_i for which $f(X_i)$ is in a bin is equal; such binning ideas are central to the theory of calibration. Then we choose the postprocessed function g by averaging observed Y values over those bins. This transforms the (population level) idea present in Theorem 15.1.1, which says to choose the post-processing conditional expectation $g(x) = \mathbb{E}[Y \mid f(X) = f(x)]$, into one implementable in finite samples, which approximately sets

$$g(x) \approx \mathbb{E}[Y \mid l(x) \leq f(X) \leq u(x)],$$

where l and u are lower and upper bounds over which to average the predictions of f .

To make the ideas concrete, assume we have a sample $(X_i, Y_i)_{i=1}^{2n}$ of size $2n$ drawn i.i.d. according to P (where we choose $2n$ for notational convenience), which we divide into samples $\{(X_i, Y_i)\}_{i=1}^n$ and $\{(X_i, Y_i)\}_{i=n+1}^{2n}$, letting $P_n^{(1)}$ denote the empirical distribution on the first sample and $P_n^{(2)}$ that on the second. We use the first to choose the binning (quantization) of f and the second to actually choose values for the binned function. Fix a number of bins $b \in \mathbb{N}$ to be chosen, for convenience assuming that b divides n . Let the indices i_1, \dots, i_n sort $f(X_i)$, so that

$$f(X_{i_1}) < f(X_{i_2}) < \dots < f(X_{i_n}),$$

and construct index partitions I_j , $j = 1, \dots, b$, by $I_j := \{i_{b(j-1)+1}, \dots, i_{bj}\}$. Here, we have assumed (essentially) without loss of generality that the predictions $f(X_i)$ are distinct with probability 1.¹ Given this partitioning of indices I_1, \dots, I_b , for $j = 1, \dots, b$ define the lower and upper bin boundaries

$$\hat{l}_j = \max_{i \in I_{j-1}} f(X_i) \quad \text{and} \quad \hat{u}_j = \max_{i \in I_j} f(X_i),$$

except that $\hat{l}_1 = 0$ and $\hat{u}_b = 1$, and define the bins

$$B_1 = [\hat{l}_1, \hat{u}_1), \quad B_2 = [\hat{l}_2, \hat{u}_2), \dots, B_b = [\hat{l}_b, \hat{u}_b]$$

to partition $[0, 1]$. These partition $[0, 1]$ evenly in the empirical probabilities of $f(X_i)$, $i = 1, \dots, n$, not evenly in the widths $\hat{u}_j - \hat{l}_j$.

To construct the recalibrated and binned version g of f , for each $x \in \mathcal{X}$, define the bin mapping

$$\text{bin}(x) := \text{the bin } j \text{ such that } f(x) \in B_j,$$

¹If this distinctness fails, we can add random dithering by letting $U_i \stackrel{\text{iid}}{\sim} \text{Uniform}[-\frac{1}{2}, \frac{1}{2}]$ and replacing the observations X_i with pairs (X_i, U_i) and $f(X_i)$ with $f_{\text{ext}}(X_i, U_i) := f(X_i) + \epsilon U_i$ for some $\epsilon > 0$. Then $L(f_{\text{ext}}) = \mathbb{E}[(Y - f(X) - \epsilon U)^2] = \mathbb{E}[(Y - f(X))^2] + \frac{\epsilon^2}{12}$ and $\ell(f_{\text{ext}}(x, u), y) \leq \ell(f(x), y) + 2\epsilon$ for all x, u, y , so that we lose little.

which implicitly depends on the first sample (X_1^n, Y_1^n) . The partitioning of $[0, 1]$ into the bins B_j also induces a partition on $\mathcal{X} = \bigcup_{j=1}^b f^{-1}(B_j)$, where elements x, x' belong to the same partition set if $\text{bin}(x) = \text{bin}(x')$. Once we have this mapping from x to the associated prediction bin, we can use the second sample (its empirical distribution) to define the binned function g by the average of the second sample distribution $P_n^{(2)}$ over those examples falling into each bin. Formally, we define g to be the piecewise constant function

$$g(x) := \mathbb{E}_{P_n^{(2)}}[Y \mid \text{bin}(X) = \text{bin}(x)], \quad (15.4.1)$$

or equivalently, for each $x \in B_j$, we have

$$\begin{aligned} g(x) &:= \mathbb{E}_{P_n^{(2)}}[Y \mid f(X) \in B_j] \\ &= \frac{1}{\sum_{i=n+1}^{2n} \mathbf{1}\{\text{bin}(X_i) = j\}} \sum_{i=n+1}^{2n} \mathbf{1}\{\text{bin}(X_i) = j\} Y_i \end{aligned}$$

where we assign $g(x)$ an arbitrary value if no X_i satisfies $f(X_i) \in B_j$ for the index $j = \text{bin}(x)$.

Informally, this function g partitions X space into regions of roughly equal (small) probability $1/b$, and for which $f(x)$ belongs to a given interval on each region. Then recalibrating f on that region changes the prediction error $(Y - f(X))^2$ little, but improves the calibration. Formally, we can show the following theorem.

Theorem 15.4.1. *Let g be the binned and recalibrated estimator (15.4.1). Assume that the number of bins b and sample size n satisfy $\frac{n}{\log n} \geq b$. Then there exists a numerical constant $c > 0$ such that for all $\delta \in (0, 1)$, with probability at least $1 - 2\exp(-c\frac{n}{b}) - \delta$,*

$$L(g) \leq L(f) + \frac{3}{b} + \frac{2b \log \frac{2b}{\delta}}{n} - \mathbb{E} \left[(\mathbb{E}[Y \mid \text{bin}(X)] - \mathbb{E}[f(X) \mid \text{bin}(X)])^2 \right]$$

and g has expected calibration error (15.2.1) at most

$$\text{ece}(g) \leq \sqrt{\frac{2b \log \frac{2b}{\delta}}{n}}.$$

JCD Comment: Put in some figures here.

The proof of Theorem 15.4.1 is long, so we defer it to Section 15.4.1. To interpret the theorem, consider the terms in it. Roughly, we see that if we choose the number of bins to be $\sqrt{n \log \frac{1}{\delta}}$, then the calibrating predictor g guarantees

$$L(g) \leq L(f) + O(1) \sqrt{\frac{\log \frac{n}{\delta}}{n}} - \mathbb{E} \left[(\mathbb{E}[Y \mid \text{bin}(X)] - \mathbb{E}[f(X) \mid \text{bin}(X)])^2 \right],$$

while the expected calibration error is of order $n^{-1/4}$, ignoring the logarithmic factors. So we improve the loss $L(f)$ by a factor involving the calibration error of f (relative to the random binning)—the less calibrated f is, the more improvement we can provide—and with a penalty tending to 0 at rate $\sqrt{\log n/n}$.

15.4.1 Proof of Theorem 15.4.1

Throughout the proof, we use the shorthands that $P(B_j) = P(f(X) \in B_j)$ and $P_n(B_j) = P_n(f(X) \in B_j) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{f(X_i) \in B_j\}$ to mean the (empirical) probability that $f(X) \in B_j$, and $P_n^{(1)}$ and $P_n^{(2)}$ denote empirical probabilities with respect to the samples (X_1^n, Y_1^n) and $(X_{n+1}^{2n}, Y_{n+1}^{2n})$, respectively. The key to the argument is to show three things:

1. With high probability, each bin B_j has the approximately correct probability $\frac{1}{2b} \leq P(B_j) \leq \frac{7}{4b}$.
2. With similarly high probability, the empirical probabilities on the second sample $P_n^{(2)}$ satisfy $\frac{1}{4b} \leq P_n^{(2)}(B_j) \leq \frac{2}{b}$.
3. Conditional on $P_n^{(2)}(B_j)$ being large enough, the expectations $\mathbb{E}_{P_n^{(2)}}[Y \mid f(X) \in B_j]$ are accurate, so that $g(x) \approx \mathbb{E}[Y \mid f(X) \in B_j]$ for x satisfying $f(x) \in B_j$.

Once we have each of these three, we can show that $L(g)$ is essentially no larger than $L(f)$, up to diminishing error terms in n , and that g itself is well-calibrated. We proceed through each step in turn, stating the results as lemmas whose proofs we provide at the end of this section.

Lemma 15.4.2. *Let $\frac{n}{\log n} \geq b$. For a numerical constant $c > 0$, we have*

$$\mathbb{P}\left(\frac{1}{2b} \leq P(B_j) \leq \frac{7}{4b} \text{ for all } j = 1, \dots, b\right) \geq 1 - 2 \exp\left(-c \frac{n}{b}\right).$$

With Lemma 15.4.2 in hand, the second step of the proof of Theorem 15.4.1 is relatively straightforward. In the lemma, conditioning on $P_n^{(1)}$ indicates conditioning on the first sample (X_1^n, Y_1^n) .

Lemma 15.4.3. *Let $\frac{n}{\log n} \geq b$. Assume the first sample $P_n^{(1)}$ is such that $\frac{1}{2b} \leq P(B_j) \leq \frac{7}{4b}$ for each selected bin B_j , $j = 1, \dots, b$. Then there exists a numerical constant $c > 0$ such that*

$$\mathbb{P}\left(\frac{1}{4b} \leq P_n^{(2)}(B_j) \leq \frac{2}{b} \mid P_n^{(1)}\right) \geq 1 - 2 \exp\left(-c \frac{n}{b}\right).$$

Lemma 15.4.4. *Let the conditions of Lemma 15.4.3 hold. Then there exists a numerical constant $c > 0$ such that for any $\delta \in (0, 1)$*

$$\mathbb{P}\left(\max_{j \leq b} \sup_{x: f(x) \in B_j} |g(x) - \mathbb{E}[Y \mid f(X) \in B_j]| \geq \sqrt{\frac{2b}{n} \log \frac{2b}{\delta}} \mid P_n^{(1)}\right) \leq 2 \exp\left(-c \frac{n}{b}\right) + \delta.$$

With the three lemmas in place, we can now expand the squared error to obtain the calibrating theorem. Recalling the population squared error $L(g) = \mathbb{E}[(Y - g(X))^2]$, let us suppose that the consequences of Lemmas 15.4.2–15.4.4 hold, so that $|g(x) - \mathbb{E}[Y \mid f(X) \in B_j]|^2 \leq \frac{2b}{n} \log \frac{2b}{\delta}$ and $P(B_j) \leq \frac{7}{4b}$ for each j . By the lemmas, these hold with probability $1 - 2 \exp(-c \frac{n}{b}) - \delta$. Define the average function values and conditional expectations

$$\bar{f}_j := \mathbb{E}[f(X) \mid f(X) \in B_j] \quad \text{and} \quad \bar{E}_j := \mathbb{E}[Y \mid f(X) \in B_j].$$

Then we have

$$L(g) = \mathbb{E}[(Y - g(X))^2] = \sum_{j=1}^b P(B_j) \mathbb{E}[(Y - \bar{E}_j + \bar{E}_j - g(X))^2 \mid f(X) \in B_j].$$

Considering the expectation term, note that $g(X)$ is constant for $f(X) \in B_j$ by construction of the binning, and so for any $x \in f^{-1}(B_j)$, we have

$$\begin{aligned} & \mathbb{E}[(Y - \bar{E}_j + \bar{E}_j - g(X))^2 \mid f(X) \in B_j] \\ &= \mathbb{E}[(Y - \mathbb{E}[Y \mid f(X) \in B_j])^2 \mid f(X) \in B_j] + (g(x) - \mathbb{E}[Y \mid f(X) \in B_j])^2 \\ &\leq \mathbb{E}[(Y - \mathbb{E}[Y \mid f(X) \in B_j])^2 \mid f(X) \in B_j] + \frac{2b}{n} \log \frac{2b}{\delta}. \end{aligned}$$

Now, using that $\mathbb{E}[Y \mid f(X) \in B_j] = \bar{E}_j$, we see that

$$\mathbb{E}[(Y - \bar{E}_j)^2 \mid f(X) \in B_j] = \mathbb{E}[(Y - \bar{f}_j)^2 \mid f(X) \in B_j] - (\bar{E}_j - \bar{f}_j)^2$$

by adding and subtracting \bar{f}_j and expanding the square. Summarizing, we have shown so far that

$$L(g) \leq \sum_{j=1}^b P(B_j) \mathbb{E}[(Y - \bar{f}_j)^2 \mid f(X) \in B_j] + \frac{2b}{n} \log \frac{2b}{\delta} - \sum_{j=1}^b P(B_j) (\bar{E}_j - \bar{f}_j)^2. \quad (15.4.2)$$

We can directly relate the first term in the expansion (15.4.2) to the expected error $\mathbb{E}[(Y - f(X))^2]$. Indeed, by expanding out the square, we have

$$\begin{aligned} & \mathbb{E}[(Y - \bar{f}_j)^2 \mid f(X) \in B_j] \\ &= \mathbb{E}[(Y - f(X) + f(X) - \bar{f}_j)^2 \mid f(X) \in B_j] \\ &= \mathbb{E}[(Y - f(X))^2 \mid f(X) \in B_j] + 2\mathbb{E}[(Y - f(X))(f(X) - \bar{f}_j) \mid f(X) \in B_j] + \text{Var}(f(X) \mid f(X) \in B_j) \\ &\leq \mathbb{E}[(Y - f(X))^2 \mid f(X) \in B_j] + 2\sqrt{\text{Var}(f(X) \mid f(X) \in B_j) + \text{Var}(f(X) \mid f(X) \in B_j)}, \end{aligned}$$

where the inequality is Cauchy-Schwarz, as $|Y - f(X)| \leq 1$. Finally, we recognize that $B_j \subset [\hat{l}_j, \hat{u}_j]$, so $\text{Var}(f(X) \mid f(X) \in B_j) \leq \frac{1}{4}(\hat{u}_j - \hat{l}_j)^2$, and thus

$$\mathbb{E}[(Y - \bar{f}_j)^2 \mid f(X) \in B_j] \leq \mathbb{E}[(Y - f(X))^2 \mid f(X) \in B_j] + \frac{5}{4}(\hat{u}_j - \hat{l}_j).$$

Substituting in the bound (15.4.2) and recognizing that $\sum_{j=1}^b P(B_j) \mathbb{E}[(Y - f(X))^2 \mid f(X) \in B_j] = \mathbb{E}[(Y - f(X))^2] = L(f)$, we have

$$L(g) \leq L(f) + \frac{5}{4} \sum_{j=1}^b P(B_j) (\hat{u}_j - \hat{l}_j) + \frac{2b}{n} \log \frac{2b}{\delta} - \sum_{j=1}^b P(B_j) (\bar{E}_j - \bar{f}_j)^2.$$

But of course, $P(B_j) \leq \frac{7}{4b}$ by the assumed conclusions of Lemma 15.4.2, and so $\sum_{j=1}^b P(B_j) (\hat{u}_j - \hat{l}_j) \leq \frac{7}{4b}$ as $\sum_{j=1}^b (\hat{u}_j - \hat{l}_j) = 1$. This gives the final inequality

$$L(g) \leq L(f) + \frac{35}{16b} + \frac{2b}{n} \log \frac{2b}{\delta} - \sum_{j=1}^b P(B_j) (\bar{E}_j - \bar{f}_j)^2,$$

proving the first claim of the theorem. The bound on calibration error is immediate because $|g(x) - \mathbb{E}[Y \mid f(X) \in B_j]|^2 \leq \frac{2b}{n} \log \frac{2b}{\delta}$ for each $x \in f^{-1}(B_j)$ with the prescribed probability, by Lemma 15.4.4.

Proof of Lemma 15.4.2 We follow the notational shorthand $P_n(A) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{f(X_i) \in A\}$. Fix a pair $0 \leq l < u \leq 1$ and define the interval $A = [l, u]$. Then Bernstein's inequality (4.1.8) shows that

$$\mathbb{P}\left(\left|\frac{1}{n}P_n(A) - P(A)\right| \geq v\right) \leq 2 \exp\left(-\frac{nv^2}{2P(A) + \frac{2}{3}v}\right)$$

for all $v \geq 0$. Partition $[0, 1]$ into intervals A_1, \dots, A_{4b} , $A_j = [l_j, u_j]$, each of probability $P(A_j) = \frac{1}{4b}$. Now, fix an index $j^* \in [b]$ and consider the (empirically constructed) bin $B_{j^*} = [\hat{l}_{j^*}, \hat{u}_{j^*}]$. Then there exist some $j, k \in \mathbb{N}$ such that

$$A_j \cup \dots \cup A_{j+k} \supset B_{j^*} \supset A_{j+1} \cup \dots \cup A_{j+k-1}.$$

We provide upper and lower bounds on k as a function of the error in $P_n(A_j)$. Suppose that for some $t > 0$, we have

$$\frac{1-t}{4b} \leq P_n(A_j) \leq \frac{1+t}{4b} \quad \text{for } j = 1, \dots, 4b. \quad (15.4.3)$$

Then

$$\frac{1+t}{4b}(k+1) \geq P_n(A_j \cup \dots \cup A_{j+k}) \geq P_n(B_{j^*}) = \frac{1}{b},$$

and similarly

$$\frac{1-t}{4b}(k-1) \leq P_n(A_{j+1} \cup \dots \cup A_{j+k}) \leq P_n(B_{j^*}) = \frac{1}{b},$$

implying the bounds

$$\frac{4}{1+t} - 1 \leq k \leq \frac{4}{t-1} + 1.$$

In particular, if $t < \frac{1}{3}$ then $3 \leq k \leq 6$, and so when the bounds (15.4.3) hold with $t = \frac{1}{3}$ we obtain

$$\frac{1}{2b} \leq \frac{k-1}{4b} = P(A_{j+1} \cup \dots \cup A_{j+k-1}) \leq P(B_{j^*}) \leq P(A_j \cup \dots \cup A_{j+k}) = \frac{k+1}{4b} \leq \frac{7}{4b}.$$

Apply Bernstein's inequality for using $t = \frac{1}{3}$, or $v = \frac{1}{12b}$, with variance bound $\sigma^2 \leq P(A_j) \leq \frac{1}{4b}$ to obtain that for each $j = 1, \dots, 4b$, we have

$$\mathbb{P}\left(|P_n(A_j) - P(A_j)| \geq \frac{1}{12b}\right) \leq 2 \exp\left(-\frac{n/(12b)^2}{2/(4b) + \frac{2}{3}\frac{1}{12b}}\right) = 2 \exp\left(-\frac{n}{80b}\right).$$

Apply a union bound to obtain the lemma once we recognize that $n/b - \log b \gtrsim n/b$ whenever $n/\log n \geq b$. \square

Proof of Lemma 15.4.3 Assume that $P(B_j) \leq \frac{7}{4b}$. Then applying Bernstein's inequality (4.1.8), and using that $\mathbf{1}\{f(X) \in B_j\}$ is a Bernoulli random variable with mean (and hence variance) at most $\frac{7}{4b}$, we have

$$\mathbb{P}\left(P_n^{(2)}(B_j) \geq \frac{2}{b}\right) \leq \exp\left(-\frac{n/(4b)^2}{\frac{7}{4b} + \frac{2}{3}\frac{1}{4b}}\right) = \exp\left(-\frac{1}{28 + 8/3} \frac{n}{b}\right) \leq \exp\left(-\frac{1}{31} \frac{n}{b}\right).$$

Similarly, we have $\mathbb{P}(P_n^{(2)}(B_j) \leq \frac{1}{4b}) \leq \exp(-\frac{1}{31} \frac{n}{b})$ as $P(B_j) \geq \frac{1}{2b}$. Applying a union bound over $j = 1, \dots, b$, then noting that $n/b - \log b \gtrsim n/b$ whenever $n/\log n \geq b$, we again obtain \square

Proof of Lemma 15.4.4 Recall that $g(x) = \mathbb{E}_{P_n^{(2)}}[Y \mid \text{bin}(X) = \text{bin}(x)]$, and note that g is constant on $x \in f^{-1}(B_j)$. Fix a bin j , and let $I(j) = \{i \in \{n+1, \dots, 2n\} \mid f(X_{n+i}) \in B_j\}$ denote the indices in the second sample for which $f(X_{n+i})$ falls in bin B_j . Then conditional on $i \in I(j)$, we have $Y_i \sim P(Y \in \cdot \mid f(X) \in B_j)$, so that

$$\mathbb{P} \left(\left| \frac{1}{|I(j)|} \sum_{i \in I(j)} Y_i - \mathbb{E}[Y \mid f(X) \in B_j] \right| \geq t \mid I(j) \right) \leq 2 \exp(-2 \text{card}(I(j))t^2)$$

by Hoeffding's inequality. Then (conditioning on the bins $\{B_j\}$ chosen using $P_n^{(1)}$, which by assumption satisfy $P(B_j) \in [\frac{1}{2b}, \frac{7}{4b}]$, we have for any fixed $x \in f^{-1}(B_j)$ that

$$\begin{aligned} & \mathbb{P} \left(\sup_{x \in f^{-1}(B_j)} |g(x) - \mathbb{E}[Y \mid f(X) \in B_j]| \geq t \mid P_n^{(1)} \right) \\ &= \sum_{I \subset [n]} \mathbb{P} \left(|g(x) - \mathbb{E}[Y \mid f(X) \in B_j]| \geq t, I(j) = I \mid P_n^{(1)} \right) \\ &\leq \mathbb{P} \left(\text{card}(I(j)) < \frac{n}{4b} \mid P_n^{(1)} \right) + \sum_{I \subset [n], \text{card}(I) \geq n/4b} \mathbb{P} \left(|g(x) - \mathbb{E}[Y \mid f(X) \in B_j]| \geq t, I(j) = I \mid P_n^{(1)} \right) \\ &\leq \mathbb{P} \left(P_n^{(2)}(B_j) < \frac{1}{4b} \right) + 2 \exp \left(-\frac{nt^2}{2b} \right), \end{aligned}$$

where the final line applies Hoeffding's inequality. Taking $t^2 = \frac{2b \log \frac{2b}{\delta}}{n}$ and applying Lemma 15.4.3 and a union bound gives Lemma 15.4.4. \square

15.5 Continuous and equivalent calibration measures

We finally return to constructing a calculus and tools with which to measure calibration, addressing the issues of discontinuity of ece that Example 15.2.2 highlights, and building to a combination of results that imply Corollary 15.2.6. In the end, we will see that for appropriate classes \mathcal{F} of predictors, several potential measures $M : \mathcal{F} \rightarrow \mathbb{R}_+$ are roughly equivalent sound and complete calibration measures, all enjoying similar continuity properties. We begin with two definitions.

Definition 15.1. A function $M : \mathcal{F} \rightarrow \mathbb{R}_+$ is a continuous calibration measure for the distribution P on $\mathcal{X} \times \mathcal{Y}$ if

- (i) it is sound and complete (15.2.2), that is, $M(f) = 0$ if and only if f is calibrated for P , and
- (ii) it is continuous with respect to the $L^1(P)$ metric on \mathcal{F} , that is, for any f , if f_n is a sequence of functions with $\mathbb{E}[\|f(X) - f_n(X)\|] \rightarrow 0$, then

$$M(f) - M(f_n) \rightarrow 0.$$

A stronger definition replaces continuity with a Lipschitz requirement.

Definition 15.2. A function $M : \mathcal{F} \rightarrow \mathbb{R}_+$ is a Lipschitz calibration measure for the distribution P on $\mathcal{X} \times \mathcal{Y}$ if it is sound and complete (Definition 15.1, part (i)), and instead of part (ii) satisfies (iii) it is Lipschitz continuous with respect to the $L^1(P)$ metric on \mathcal{F} , that is, for some $C < \infty$

$$|M(f_0) - M(f_1)| \leq C \cdot \mathbb{E}_P[\|f_0(X) - f_1(X)\|]$$

for all $f_0, f_1 \in \mathcal{F}$.

If conditions (i) and (ii) (respectively (iii)) hold for all P in a collection of distributions \mathcal{P} on $\mathcal{X} \times \mathcal{Y}$, we will say that M is a continuous (respectively, Lipschitz) calibration measure for \mathcal{P} .

The desiderata (ii) and (iii) are matters of taste; the central idea is that some type of continuity is essential for efficient modeling, estimation, and analysis. We leave the norm $\|\cdot\|$ implicit in the definition, and we typically omit the distribution P from the calibration metric as it is clear from context. The two parts of Definition 15.2 admit many possible calibration measures. We consider two types of measures, which are (almost) dual to one another, as examples. Both use a variational representation, where in one we essentially look for the “closest” function that is calibrated, while in the other, we investigate the ease with which we can (quantitatively) certify that a predictor f is uncalibrated.

A key concept will be the equivalence of calibration measures, where we target a quantitative equivalence. To define this, let $0 < \alpha, \beta < \infty$. Then we say that two candidate calibration measures M_0 and M_1 on $\mathcal{F} \subset \mathcal{X} \rightarrow \mathbb{R}^k$ are (α, β) -equivalent if there exist constants c_0, c_1 (which may depend on \mathcal{Y}) such that

$$M_0(f) \leq c_0 [M_1(f) + M_1(f)^\alpha] \quad \text{and} \quad M_1(f) \leq c_1 [M_0(f) + M_0(f)^\beta]. \quad (15.5.1)$$

Then in a strong sense, $M_0(f) \rightarrow 0$ if and only if $M_1(f) \rightarrow 0$.

15.5.1 Calibration measures

We revisit the potential calibration measures in Section 15.2.2 here to recapitulate definitions, providing initial results on their soundness and completeness. We focus on the distance to calibration (15.2.3) and relative calibration errors (15.2.4), as the partitioned calibration error (15.2.6) we use more as a proof device.

Distances to calibration. Recall the *distance to calibration* (15.2.3), which for $\mathcal{C}(P) = \{g : \mathcal{X} \rightarrow \mathbb{R}^k \mid \mathbb{E}_P[Y \mid g(X)] = g(X)\}$ (where the defining equality holds with P -probability 1 over X) has definition $d_{\text{cal}}(f) := \inf_g \{\mathbb{E}[\|g(X) - f(X)\|] \text{ s.t. } g \in \mathcal{C}(P)\}$. The measure (15.2.3) is, after appropriate normalization, the *largest* Lipschitz measure of calibration: if M is any Lipschitz calibration measure (with constant $C = 1$ in Definition 15.2 part (iii)), then taking a perfectly calibrated g with $\text{ece}(g) = 0$, we necessarily have $M(g) = 0$. Then for any f we have $M(f) = M(f) - M(g) \leq \mathbb{E}[\|f(X) - g(X)\|]$, and taking an infimum over such g guarantees

$$M(f) \leq d_{\text{cal}}(f).$$

The second related quantity, which sometimes admits cleaner properties for analysis, is the *penalized calibration distance*, which we define as

$$p_{\text{cal}}(f) := \inf_g \{\mathbb{E}[\|f(X) - g(X)\|] + \mathbb{E}[\|\mathbb{E}[Y \mid g(X)] - g(X)\|]\}. \quad (15.5.2)$$

These quantities are strongly related, and in the sequel (see Corollary 15.5.8), we show that

$$p_{\text{cal}}(f) \leq d_{\text{cal}}(f) \leq p_{\text{cal}}(f) + C_{\mathcal{Y}} \sqrt{p_{\text{cal}}(f)},$$

where $C_{\mathcal{Y}}$ is a constant depending only on the set \mathcal{Y} whenever \mathcal{Y} has finite diameter.

To build intuition for the definition (15.5.2), consider the two quantities. The first measures the usual L^1 distance between the function f and a putative alternative g . The second is the expected calibration error of g . By restricting the infimum in definition (15.2.3) to functions g with $\text{ece}(g) = 0$, we simply have the L^1 distance to the nearest calibrated function; as is, the additional term in (15.5.2) allows trading between the distance to a calibrated function and the actual calibration error. We also have the following proposition.

Proposition 15.5.1. *The functions d_{cal} and p_{cal} are Lipschitz calibration measures.*

Proof If f is calibrated, then $p_{\text{cal}}(f) = d_{\text{cal}}(f) = 0$ immediately. Conversely, if $p_{\text{cal}}(f) = 0$, there exists a sequence of functions g_n satisfying $\mathbb{E}[\|f(X) - g_n(X)\|] \rightarrow 0$, as each term in the definition (15.5.2) is nonnegative. Additionally, we must have that $\text{ece}(g_n) = \mathbb{E}[\|\mathbb{E}[Y | g_n(X)] - g_n(X)\|] \rightarrow 0$. Applying Lemma 15.2.1 we have $0 \geq \liminf_n \text{ece}(g_n) \geq \text{ece}(f)$. If $d_{\text{cal}}(f) = 0$, then there exists a sequence of functions g_n with $\text{ece}(g_n) = 0$ and $\mathbb{E}[\|f(X) - g_n(X)\|] \rightarrow 0$. Again, the lower semicontinuity of ece from Lemma 15.2.1 gives $0 = \liminf_n \text{ece}(g_n) \geq \text{ece}(f)$.

To see that p_{cal} is Lipschitz in f , let $f_0, f_1 : \mathcal{X} \rightarrow \mathbb{R}^k$, and let g_0, g_1 be within $\epsilon > 0$ of achieving the infima in definition (15.5.2) for f_0 and f_1 , respectively. Then

$$\begin{aligned} p_{\text{cal}}(f_0) - p_{\text{cal}}(f_1) &\leq \inf_g \{ \mathbb{E}[\|f_0(X) - g(X)\|] + \mathbb{E}[\|\mathbb{E}[Y | g(X)] - g(X)\|] \\ &\quad - \mathbb{E}[\|f_1(X) - g_1(X)\|] + \mathbb{E}[\|\mathbb{E}[Y | g_1(X)] - g_1(X)\|] + \epsilon \\ &\leq \mathbb{E}[\|f_0(X) - g_1(X)\|] - \mathbb{E}[\|f_1(X) - g_1(X)\|] + \epsilon \\ &\leq \mathbb{E}[\|f_0(X) - f_1(X)\|] + \epsilon. \end{aligned}$$

Take $\epsilon \downarrow 0$. The lower inequality is similar, as is the proof for d_{cal} . □

Weak calibration. The calibration error (15.2.4) relative to a class \mathcal{W} ,

$$\text{CE}(f, \mathcal{W}) := \sup_{w \in \mathcal{W}} \mathbb{E}[\langle w(f(X)), Y - f(X) \rangle]$$

admits similar properties, as it also satisfies our desiderata for a calibration measure. In particular, if we take \mathcal{W} to be the class $\mathcal{W}_{\|\cdot\|}$ of bounded Lipschitz witness functions (15.2.5), we have the next two propositions.

Proposition 15.5.2. *Let \mathcal{F} consist of functions with $\mathbb{E}[\|f(X)\|] < \infty$ and assume $\mathbb{E}[\|Y\|] < \infty$. Then $\text{CE}(\cdot, \mathcal{W}_{\|\cdot\|})$ is a continuous calibration measure over \mathcal{F} .*

Because continuity is such a weak requirement, the proof of this result relies on measure theoretic results, so we defer it to Section 15.6.2.

When we assume the collection \mathcal{F} consists of bounded functions and \mathcal{Y} itself is bounded, we can give a stronger guarantee for the weak calibration, and we no longer need to rely on careful arguments considering the order of various limits.

Proposition 15.5.3. *Assume that $\text{diam}(\mathcal{Y})$ is finite and that \mathcal{F} is a collection of bounded functions $\mathcal{X} \rightarrow \mathbb{R}^k$. Then $\text{CE}(\cdot, \mathcal{W}_{\|\cdot\|})$ is a Lipschitz calibration measure over \mathcal{F} .*

Proof Let $\mathcal{W} = \mathcal{W}_{\|\cdot\|}$ for shorthand. That $\text{CE}(f, \mathcal{W}) = 0$ when f is calibrated is immediate, as by definition of conditional expectation we have

$$\mathbb{E}[\langle w(f(X)), Y - f(X) \rangle] = \mathbb{E}[\langle w(f(X)), \mathbb{E}[Y | f(X)] - f(X) \rangle] = 0.$$

To obtain the converse that $\text{CE}(f, \mathcal{W}) = 0$ implies f is calibrated, we require an intermediate lemma, which leverages the density of Lipschitz functions in L^p spaces. As was the case for the lower semi-continuity lemma 15.2.1 central to the proof of the converse in Proposition 15.5.1, this lemma requires measure-theoretic approximation arguments, so we defer its proof to Section 15.6.3.

Lemma 15.5.4. *Let $S \in \mathbb{R}^k$ be a random variable and $\mathbb{E}[\|g(S)\|] < \infty$. If $\mathbb{E}[\langle w(S), g(S) \rangle] = 0$ for all bounded and 1-Lipschitz functions w , then $g(S) = 0$ with probability 1.*

The converse is now trivial: let $S = f(X)$, and note that $\text{CE}(f, \mathcal{W}) = \sup_{w \in \mathcal{W}} \mathbb{E}[\langle w(S), \mathbb{E}[Y | S] - S \rangle]$, and take $g(S) = \mathbb{E}[Y | S] - S$ in Lemma 15.5.4.

To see that CE is Lipschitz, let $w_0 \in \mathcal{W}$ be such that $\text{CE}(f_0, \mathcal{W}) \geq \mathbb{E}[\langle w_0(f_0(X)), Y - f_0(X) \rangle] - \epsilon$, and let $C < \infty$ satisfy $C \geq \sup_{y \in \mathcal{Y}, x \in \mathcal{X}, f \in \mathcal{F}} \|y - f(x)\|$. Then

$$\begin{aligned} \text{CE}(f_0, \mathcal{W}) - \text{CE}(f_1, \mathcal{W}) &\leq \mathbb{E}[\langle w_0(f_0(X)), Y - f_0(X) \rangle] - \mathbb{E}[\langle w_0(f_1(X)), Y - f_1(X) \rangle] + \epsilon \\ &\leq \mathbb{E}[\langle w_0(f_0(X)) - w_0(f_1(X)), Y - f_0(X) \rangle] + \mathbb{E}[\langle w_0(f_1(X)), f_1(X) - f_0(X) \rangle] + \epsilon \\ &\leq C \mathbb{E}[\|w_0(f_0(X)) - w_0(f_1(X))\|_*] + \mathbb{E}[\|f_1(X) - f_0(X)\|] + \epsilon \\ &\leq (1 + C) \mathbb{E}[\|f_1(X) - f_0(X)\|] + \epsilon. \end{aligned}$$

Repeating the same argument, *mutatis mutandis*, for the lower bound gives the Lipschitz continuity as desired. \square

The family of weak calibration measures $\text{CE}(f, \mathcal{W})$ as we vary the collection of potential witness functions \mathcal{W} yields a variety of behaviors. Different choices of \mathcal{W} can give different continuous calibration measures, where we may modify Definition 15.1 part (ii) to other notions of continuity, such as Lipschitzness with respect to $L^2(P)$ norms. We explore a few of these in the exercises at the end of the chapter.

15.5.2 Equivalent calibration measures

That all three measures $d_{\text{cal}}(f)$, $p_{\text{cal}}(f)$, $\text{CE}(f, \mathcal{W}_{\|\cdot\|})$ are Lipschitz calibration measures when the label space \mathcal{Y} is bounded suggests deeper relationships between these and other notions of calibration, such as the equivalence (15.5.1). We elucidate this here, showing that each of the measures d_{cal} , p_{cal} , and CE are equivalent. Indeed, the main consequence of the results in this chapter is that this equivalence holds for multiclass classification.

Theorem 15.5.5. *Let $\mathcal{Y} = \{e_1, \dots, e_k\}$ and $\mathcal{W}_{\|\cdot\|}$ be the collection (15.2.5) of bounded Lipschitz functions for a norm $\|\cdot\|$ on \mathbb{R}^k . Then d_{cal} , p_{cal} , and $\text{CE}(\cdot, \mathcal{W}_{\|\cdot\|})$ are each $(\frac{1}{2}, \frac{1}{2})$ -equivalent. Moreover, this equivalence is sharp, in that they are not (α, β) -equivalent for any $\alpha, \beta > \frac{1}{2}$.*

The theorem follows as a compilation of the other results in this section. Along the way to demonstrating this theorem, we introduce a few alternative measures of calibration we use as stepping stones toward our final results. While many of our derivations will apply for general sets \mathcal{Y} , in some cases we will restrict to multiclass classification problems, so that $\mathcal{Y} = \{e_1, \dots, e_k\} \subset \mathbb{R}^k$ are the k standard basis vectors. We present two main results: the first, Theorem 15.5.6, shows an equivalence (up to a square root) between the penalized calibration distance (15.5.2) and the partitioned calibration error (15.2.6). As a corollary of this result, we obtain the equivalence of the distance to calibration (15.2.3) and penalized distance to calibration (15.5.2). The second main result, Theorem 15.5.9, gives a similar equivalence between the penalized distance (15.5.2) and the calibration error relative to Lipschitz functions (15.2.4). Throughout, to make the calculations cleaner and more transparent, we restrict our functions to make predictions in $\mathcal{M} = \text{conv}(\mathcal{Y})$.

Partition-based calibration measures and lifting to random variables

It is easier to work directly in the space of predictions $f(X) \in \mathbb{R}^k$ rather than addressing the underlying space \mathcal{X} . To that end, let $S = f(X)$ be the random vector (use the mnemonic that S is for “scores”) induced by $f(X)$ and taking values in $\text{Conv}(\mathcal{Y})$, which has a joint distribution (S, Y) with the label Y . Then, for example, the expected calibration error of f is simply

$$\text{ece}(f) = \mathbb{E}[\|\mathbb{E}[Y | S] - S\|].$$

Once we work exclusively in the space of random scores $S = f(X)$, we may define alternative distances to calibration in analogy with the (penalized) distances to calibration, which will allow us to more easily relate distances to the partitioned error (15.2.6). Thus, we define

$$d_{\text{cal},\text{low}}(f) := \inf_V \{\mathbb{E}[\|S - V\|] \text{ s.t. } \mathbb{E}[Y | V] = V\} \quad (15.5.3a)$$

and

$$p_{\text{cal},\text{low}}(f) := \inf_V \{\mathbb{E}[\|S - V\|] + \mathbb{E}[\|\mathbb{E}[Y | V] - V\|]\}, \quad (15.5.3b)$$

where the infimum are over all random variables V taking values in $\text{Conv}(\mathcal{Y})$, which can have arbitrary distribution with (S, Y) (but do not modify the joint (S, Y)), and in case (15.5.3a) are calibrated. This formulation is convenient in that we can represent it as a convex optimization problem, allowing us to bring the tools of duality to bear on it, though we defer this temporarily. By considering $V = g(X)$ for functions $g : \mathcal{X} \rightarrow \text{Conv}(\mathcal{Y})$, we immediately see that $p_{\text{cal}}(f) \geq p_{\text{cal},\text{low}}(f)$. We can also consider upper distances

$$d_{\text{cal},\text{up}}(f) := \inf_g \{\mathbb{E}[\|S - g(S)\|] \text{ s.t. } \mathbb{E}[Y | g(S)] = g(S)\}$$

and

$$p_{\text{cal},\text{up}}(f) := \inf_{g: \mathbb{R}^k \rightarrow \text{Conv}(\mathcal{Y})} \{\mathbb{E}[\|S - g(S)\|] + \mathbb{E}[\|\mathbb{E}[Y | g(S)] - g(S)\|]\},$$

which restrict the definitions (15.2.3) and (15.5.2) to compositions. We therefore have the inequalities

$$d_{\text{cal},\text{low}}(f) \leq d_{\text{cal}}(f) \leq d_{\text{cal},\text{up}}(f) \quad \text{and} \quad p_{\text{cal},\text{low}}(f) \leq p_{\text{cal}}(f) \leq p_{\text{cal},\text{up}}(f). \quad (15.5.4)$$

The partitioned calibration error (15.2.6) allows us to provide a bound relating the calibration error and the lower and upper calibration errors. To state the theorem, we make a normalization with $\|\cdot\|$, assuming without loss of generality that $\|\cdot\|_\infty \leq \|\cdot\|$.

Theorem 15.5.6. *Let $\mathcal{Y} \subset \mathbb{R}^k$ have finite diameter $\text{diam}(\mathcal{Y})$ in the norm $\|\cdot\|$. Let $S = f(X) \in \mathbb{R}^k$. Then for all $\varepsilon > 0$,*

$$\begin{aligned} p_{\text{cal},\text{up}}(f) \leq d_{\text{cal},\text{up}}(f) \leq \text{pce}(S) &\leq \left(1 + \frac{2k \text{diam}(\mathcal{Y})}{\varepsilon}\right) p_{\text{cal},\text{low}}(f) + \|\mathbf{1}_k\|_* \varepsilon \\ &\leq \left(1 + \frac{2k \text{diam}(\mathcal{Y})}{\varepsilon}\right) d_{\text{cal},\text{low}}(f) + \|\mathbf{1}_k\|_* \varepsilon. \end{aligned}$$

While the first inequality in Theorem 15.5.6 is relatively straightforward to prove, the second requires substantially more care, so we defer the proof of the theorem to Section 15.6.4.

We record a few corollaries, one consequence of which is to show that the partitioned calibration error (15.2.6) is at least a calibration measure in the sense of Definition 15.2.(i). Theorem 15.5.6 also shows that the penalized calibration distance $p_{\text{cal}}(f)$ is equivalent, up to taking a square root, to the upper and lower calibration “distances”. In each corollary, we let $C_k = \|\mathbf{1}_k\|_*$ for shorthand.

Corollary 15.5.7. *Let the conditions of Theorem 15.5.6 hold. Then*

$$p_{\text{cal},\text{low}}(f) \leq p_{\text{cal}}(f) \leq p_{\text{cal},\text{low}}(f) + 2\sqrt{C_k k \text{diam}(\mathcal{Y})} \sqrt{p_{\text{cal},\text{low}}(f)}$$

and

$$d_{\text{cal},\text{low}}(f) \leq d_{\text{cal}}(f) \leq d_{\text{cal},\text{low}}(f) + 2\sqrt{C_k k \text{diam}(\mathcal{Y})} \sqrt{d_{\text{cal},\text{low}}(f)}.$$

Proof The first lower bound is immediate (recall the naive inequalities (15.5.4)). Now set $\varepsilon = \sqrt{2k \text{diam}(\mathcal{Y}) p_{\text{cal},\text{low}}(f) / C_k}$ in Theorem 15.5.6, and recognize that $p_{\text{cal},\text{low}}(f) \leq p_{\text{cal},\text{up}}(f)$. \square

We also obtain an approximate equivalence between the calibration distance d_{cal} and penalized calibration distance p_{cal} from definitions (15.2.3) and (15.5.2).

Corollary 15.5.8. *Let the conditions of Theorem 15.5.6 hold. Then*

$$p_{\text{cal}}(f) \leq d_{\text{cal}}(f) \leq p_{\text{cal}}(f) + 2\sqrt{C_k k \text{diam}(\mathcal{Y})} \sqrt{p_{\text{cal}}(f)}.$$

Proof The first inequality is immediate by definition. For the second, note (see Lemma 15.6.4 in the proof of Theorem 15.5.6 in Section 15.6.4) that $p_{\text{cal},\text{low}}(f) \leq \text{pce}(S)$ for $S = f(X)$. Then apply Theorem 15.5.6 with $\varepsilon = \sqrt{2k \text{diam}(\mathcal{Y}) p_{\text{cal},\text{low}}(f) / C_k}$ as in Corollary 15.5.7, and recognize that $p_{\text{cal},\text{low}} \leq p_{\text{cal}}$. \square

Let us instantiate the theorem and its corollaries in a few special cases. If we make binary predictions with $\mathcal{Y} = \{0, 1\}$, then $C_k = k = \text{diam}(\mathcal{Y}) = 1$, and Theorem 15.5.6 implies that

$$p_{\text{cal},\text{low}}(f) \leq p_{\text{cal}}(f) \leq p_{\text{cal},\text{low}}(f) + 2\sqrt{p_{\text{cal},\text{low}}(f)}.$$

For k -class multiclass classification, where we identify $\mathcal{Y} = \{e_1, \dots, e_k\}$ with the k standard basis vectors, we have the bounds

$$p_{\text{cal},\text{low}}(f) \leq p_{\text{cal}}(f) \leq p_{\text{cal},\text{low}}(f) + 2\sqrt{k p_{\text{cal},\text{low}}(f)},$$

so long as we measure calibration errors with respect to the ℓ_1 -norm, that is, $\|y - f(x)\|_1$, because $\text{diam}(\mathcal{Y}) \leq 1$ and $C_k = \|\mathbf{1}\|_\infty = 1$.

JCD Comment: Remark on sharpness here.

The equivalence between calibration error and the calibration distance

We can rewrite the calibration error $\text{CE}(S, \mathcal{A})$ relative to partitions in the definition (15.2.6) as the supremum over a collection $\mathcal{W}_{\mathcal{A}}$ of functions of the form $w(s) = v \mathbf{1}\{s \in A\}$, where $\|v\|_* \leq 1$, so that $\text{CE}(S, \mathcal{W}_{\mathcal{A}}) = \sup_{w \in \mathcal{W}_{\mathcal{A}}} \mathbb{E}[\langle w(S), Y - S \rangle] = \sum_{A \in \mathcal{A}} \mathbb{E}[\|\mathbb{E}[Y | S] - S\|]$. Relaxing this supremum, and removing the infimum over partitions, we might expect a similar relationship to Theorem 15.5.6 to hold. Via a duality argument that the definition (15.5.3) of the lower calibration error as an infimum over joint distributions makes possible, we can directly relate the measures.

Theorem 15.5.9. *Let $\mathcal{Y} \subset \mathbb{R}^k$ have finite diameter in the norm $\|\cdot\|$ and $\mathcal{W}_{\|\cdot\|}$ be the collection (15.2.5) of bounded Lipschitz functions. Then*

$$\text{CE}(f, \mathcal{W}_{\|\cdot\|}) \leq (1 + \text{diam}(\mathcal{Y})) \cdot p_{\text{cal}, \text{low}}(f).$$

Conversely, let $\mathcal{Y} = \{e_1, \dots, e_k\}$ and define $C_k := \|\mathbf{1}_k\|_ \max\{1, \text{diam}(\mathcal{Y})\}$. Then*

$$d_{\text{cal}, \text{low}}(f) \leq C_k \cdot \text{CE}(f, \mathcal{W}_{\|\cdot\|}).$$

This proof, while nontrivial, is more elementary than the others in this chapter, so we present it here. Before giving it, however, we give a few corollaries that give a fuller picture of the relationships between the different calibration measures we have developed. These show how, for the case of k -class multiclass classification where we identify $\mathcal{Y} = \{e_1, \dots, e_k\}$ with the standard basis vectors, the distance to calibration (15.2.3) and penalized calibration (15.5.2) provide essentially equivalent measures of calibration error, and that these in turn are equivalent to the calibration error with respect to the collection of bounded Lipschitz functions.

We first give a corollary for the penalized calibration (15.5.2).

Corollary 15.5.10. *Let $\mathcal{Y} = \{e_1, \dots, e_k\}$ and $\|\cdot\| = \|\cdot\|_1$ be the ℓ_1 -norm. Then for any $f : \mathcal{X} \rightarrow \text{Conv}(\mathcal{Y})$, we have*

$$\frac{1}{2} \text{CE}(f, \mathcal{W}_{\|\cdot\|}) \leq p_{\text{cal}}(f) \leq \text{CE}(f, \mathcal{W}_{\|\cdot\|}) + 2\sqrt{k \text{CE}(f, \mathcal{W}_{\|\cdot\|})}.$$

Proof Let $\mathcal{W} = \mathcal{W}_{\|\cdot\|}$ for shorthand. Theorem 15.5.9 gives $\text{CE}(f, \mathcal{W}) \leq 2p_{\text{cal}, \text{low}}(f)$, and $p_{\text{cal}, \text{low}}(f) \leq p_{\text{cal}}(f)$, giving the lower bound. For the upper bound, Corollary 15.5.7 gives $p_{\text{cal}}(f) \leq p_{\text{cal}, \text{low}}(f) + 2\sqrt{k} \sqrt{p_{\text{cal}, \text{low}}(f)}$, then using that $p_{\text{cal}, \text{low}}(f) \leq d_{\text{cal}, \text{low}}(f)$ and the second part of Theorem 15.5.9 gives the corollary. \square

The same argument implies the following analogue for the distance to calibration (15.2.3).

Corollary 15.5.11. *Let $\mathcal{Y} = \{e_1, \dots, e_k\}$ and $\|\cdot\| = \|\cdot\|_1$ be the ℓ_1 -norm. Then for any $f : \mathcal{X} \rightarrow \text{Conv}(\mathcal{Y})$, we have*

$$\frac{1}{2} \text{CE}(f, \mathcal{W}_{\|\cdot\|}) \leq d_{\text{cal}}(f) \leq \text{CE}(f, \mathcal{W}_{\|\cdot\|}) + 2\sqrt{k \text{CE}(f, \mathcal{W}_{\|\cdot\|})}.$$

Proof of Theorem 15.5.9

The proof of the upper bound is fairly straightforward. For any $w \in \mathcal{W}_{\|\cdot\|}$, we have

$$\begin{aligned} \mathbb{E}[\langle w(S), Y - S \rangle] &= \mathbb{E}[\langle w(S), V - S \rangle] + \mathbb{E}[\langle w(S) - w(V), Y - V \rangle] + \mathbb{E}[\langle w(V), Y - V \rangle] \\ &\leq \mathbb{E}[\|V - S\|] + \text{diam}(\mathcal{Y}) \mathbb{E}[\|V - S\|] + \mathbb{E}[\langle w(V), \mathbb{E}[Y | V] - V \rangle] \\ &\leq (1 + \text{diam}(\mathcal{Y})) \mathbb{E}[\|V - S\|] + \mathbb{E}[\|\mathbb{E}[Y | V] - V\|]. \end{aligned}$$

To prove the converse requires more; we present most of the argument for an arbitrary discrete space \mathcal{Y} and specialize to the multiclass setting only at the end. The starting point is to reduce the problem to a discrete problem over probability mass functions rather than general distributions, as then it is much easier to apply the standard tools of convex duality. Consider the value

$$d_{\text{cal,low}}(S) = \inf_V \{ \mathbb{E}[\|S - V\|] \text{ s.t. } \mathbb{E}[Y | V] = V \}.$$

Let $b \in \mathbb{N}$ and \mathcal{S}_b be a (minimal) $1/b$ covering $\{s_1, \dots, s_N\}$ of $\text{Conv}(\mathcal{Y})$, and define S_b to be the projection of S to the nearest s_i . Then $\|S - V\| = \|S_b - V\| \pm \frac{1}{b}$, and

$$d_{\text{cal,low}}(S) = \inf_V \{ \mathbb{E}[\|S_b - V\|] \text{ s.t. } \mathbb{E}[Y | V] = V \} \pm \frac{1}{b}.$$

Now, if we replace the infimum over arbitrary joint distributions of (S_b, Y, V) leaving the marginal (S_b, Y) unchanged (with V calibrated) with an infimum over only discrete distributions on V , we have

$$d_{\text{cal,low}}(S) \leq \inf_{V \text{ finitely supported}} \{ \mathbb{E}[\|S_b - V\|] \text{ s.t. } \mathbb{E}[Y | V] = V \} + \frac{1}{b}. \quad (15.5.5)$$

Notably, the infimum is non-empty, as we can always choose $V = Y$.

With the problem (15.5.5) in hand, we can write a finite dimensional optimization problem whose optimal value is the discretized infimum on the right side. Without loss of generality assuming that S is finitely supported, we let $p_{sy} = \mathbb{P}(S = s, Y = y)$ be the probability mass function of (S, Y) . Then introducing the joint distribution Q with p.m.f. $q_{syv} = Q(S = s, Y = y, V = v)$, the infimum (15.5.5) has the constraint that $\sum_v q_{syv} = p_{sy}$. Then $\mathbb{E}[\|S - V\|] = \sum_{s,y,v} q_{syv} \|s - v\|$ and the calibration constraint $\mathbb{E}[Y | V] = V$ is equivalent to the equality constraint that $\sum_{s,y} q_{syv}(y - v) = 0$ for each v . This yields the convex optimization problem

$$\begin{aligned} & \text{minimize} && \sum_{s,y,v} q_{syv} \|s - v\| \\ & \text{subject to} && \sum_v q_{syv} = p_{sy}, \quad q \succeq 0, \quad \sum_{s,y} q_{syv}(y - v) = 0 \text{ for all } v \end{aligned} \quad (15.5.6)$$

in the variable q . We take the dual of this problem. Taking Lagrange multipliers λ_{sy} for each equality constraint that $\sum_v q_{syv} = p_{sy}$, $\theta_{syv} \geq 0$ for the nonnegativity constraints on q , and $\beta_v \in \mathbb{R}^k$ for each equality constraint that $0 = \sum_{s,y} q_{syv}(y - v)$, we have Lagrangian

$$\begin{aligned} \mathcal{L}(q, z, \lambda, \theta, \beta) &= \sum_{s,y,v} q_{syv} \|s - v\| + \sum_{s,y,v} q_{syv} \beta_v^T (y - v) - \sum_{s,y} \lambda_{sy} \left(\sum_v q_{syv} - p_{sy} \right) - \langle \theta, q \rangle. \end{aligned}$$

Taking an infimum over q , we see that unless

$$\|s - v\| + \beta_v^T (y - v) - \lambda_{sy} - \theta_{syv} = 0$$

for each triple (s, y, v) , we have $\inf_q \mathcal{L}(q, \lambda, \theta, \beta) = -\infty$. The equality in the preceding display is equivalent to $\|s - v\| + \beta_v^T (y - v) \geq \lambda_{sy}$, so that eliminating $\theta \succeq 0$ variables, we have the dual

$$\begin{aligned} & \text{maximize} && \sum_{s,y} \lambda_{sy} p_{sy} \\ & \text{subject to} && \lambda_{sy} \leq \|s - v\| + \beta_v^T (y - v), \quad \text{all } s, y, v \end{aligned}$$

to problem (15.5.6). Equivalently, recognizing that at the optimum we must saturate the constraints on λ via $\lambda_{sy} = \min_v \{\|s - v\| + \beta_v^T(y - v)\}$, we have

$$\text{maximize} \quad \sum_{s,y} p_{sy} \min_v \{\|s - v\| + \beta_v^T(y - v)\} \quad (15.5.7)$$

in the variables β_v , and strong duality obtains.

The dual problem (15.5.7) is the key to the final step in the proof. To make the functional notation clearer, let us fix any collection of vectors β_v and define $\lambda_y(s) = \min_v \{\|s - v\| + \beta_v^T(y - v)\}$ for each $y \in \mathcal{Y}$. If we can exhibit a C -Lipschitz function $s \mapsto w(s)$ that satisfies

$$\langle w(s), y - s \rangle \geq \lambda_y(s) \quad (15.5.8)$$

for each $y \in \mathcal{Y}$ and $\|w(s)\|_* \leq C$, we will evidently have shown that

$$\sup_{w \in \mathcal{W}_{\|\cdot\|}} \mathbb{E}[\langle w(S), Y - S \rangle] \geq \frac{1}{C} d_{\text{cal,low}}(S),$$

by the dual formulation (15.5.7).

The functions λ_y are each 1-Lipschitz with respect to $\|\cdot\|$, as

$$\begin{aligned} \lambda_y(s) - \lambda_y(s') &\geq \min_v \{\|s - v\| + \beta_v^T(y - v) - \|s' - v\| - \beta_v^T(y - v)\} \\ &= \min_v \{\|s - v\| - \|s' - v\|\} \geq -\|s - s'\|, \end{aligned}$$

and similarly

$$\lambda_y(s) - \lambda_y(s') \leq \max_v \{\|s - v\| + \beta_v^T(y - v) - \|s' - v\| - \beta_v^T(y - v)\} \leq \|s - s'\|$$

by the triangle inequality. Here, we specialize to the particular multiclass classification case in which the set $\mathcal{Y} = \{e_1, \dots, e_k\}$ consists of extreme points of the probability simplex, so that $s \in \text{Conv}(\mathcal{Y})$ means that $\langle \mathbf{1}, s \rangle = 1$ and $s \succeq 0$. Abusing notation slightly, let $\lambda_i = \lambda_{e_i}$ for $i = 1, \dots, k$. Then define the function

$$w(s) := \begin{bmatrix} \lambda_1(s) \\ \vdots \\ \lambda_k(s) \end{bmatrix}.$$

By inspection, we have

$$\|w(s) - w(s')\|_* \leq \| \|s - s'\| \mathbf{1} \|_* = \|\mathbf{1}\|_* \|s - s'\|.$$

Additionally, because $\lambda_i(s) \leq \|s - e_i\|$ (take $v = e_i$ in the definition of λ_i), we have $\|w(s)\|_* \leq \|\mathbf{1}\|_* \text{diam}(\mathcal{Y})$. Finally, we have

$$\begin{aligned} \langle w(s), e_i - s \rangle &= (1 - s_i) \lambda_i(s) - \sum_{j \neq i} s_j \lambda_j(s) \\ &\geq (1 - s_i) \lambda_i(s) - \sum_{j \neq i} s_j \langle \beta_s, e_j - s \rangle \end{aligned}$$

because $\lambda_j(s) \leq \langle \beta_s, e_j - s \rangle$ by taking $v = s$ in the definition of λ_j . Adding and subtracting $s_i \langle \beta_s, e_i - s \rangle$, we obtain

$$\begin{aligned} \langle w(s), e_i - s \rangle &\geq (1 - s_i) \lambda_i(s) - \sum_{j=1}^k s_j \langle \beta_s, e_j - s \rangle + s_i \langle \beta_s, e_i - s \rangle \\ &= (1 - s_i) \lambda_i(s) + s_i \langle \beta_s, e_i - s \rangle \geq \lambda_i(s), \end{aligned}$$

because $s \succeq 0$ and $\langle \beta_s, e_i - s \rangle \geq \lambda_i(s)$. This is the desired inequality (15.5.8).

15.6 Deferred technical proofs

JCD Comment: Put these in the appendices

Several of the proofs in this chapter rely on standard results from analysis and measure theory; we give these as base lemmas, as any book on graduate level real analysis (implicitly) contains them (see, e.g., Tao [176, Chapters 1.3 and 1.13] or Royden [163]).

Lemma 15.6.1 (Egorov's theorem). *Let $f_n \rightarrow f$ in $L^p(P)$ for some $p \geq 1$. Then for each $\epsilon > 0$, there exists a set A of measure at least $P(A) \geq 1 - \epsilon$ such that $f_n \rightarrow f$ uniformly on A .*

Lemma 15.6.2 (Monotone convergence). *Let $f_n : \mathcal{X} \rightarrow \mathbb{R}_+$ be a monotone increasing sequence of functions and $f(x) = \lim_n f_n(x)$ (which may be infinite). Then $\int f(x) d\mu(x) = \lim_n \int f_n(x) d\mu(x)$ for any measure μ .*

Lemma 15.6.3 (Density of Lipschitz functions). *Let $\mathcal{C}_c^{\text{Lip}}$ be the collection of compactly supported Lipschitz functions on \mathbb{R}^k and P a probability distribution on \mathbb{R}^k . Then $\mathcal{C}_c^{\text{Lip}}$ is dense in $L^p(P)$, that is, for each $\epsilon > 0$ and f with $\mathbb{E}_P[|f(X)|^p] < \infty$, there exists $g \in \mathcal{C}_c^{\text{Lip}}$ with $\mathbb{E}_P[|g(X) - f(X)|^p]^{1/p} \leq \epsilon$.*

15.6.1 Proof of Lemma 15.2.1

Let \mathcal{W}_k be the collection of k -Lipschitz functions w with $\|w(s)\|_* \leq 1$ for all s , and let \mathcal{W} denote the collection of measurable functions with $\|w(s)\|_* \leq 1$ for all s . Recall the definition $\text{CE}(g, \mathcal{W}_k) = \sup_{w \in \mathcal{W}_k} \mathbb{E}[\langle w(g(X)), Y - g(X) \rangle]$. Then if $f_n \rightarrow f$ in $L^1(P)$, by Egorov's theorem (Lemma 15.6.1), for each $\epsilon > 0$ there exists a set A with $P(A) \geq 1 - \epsilon$ and $f_n \rightarrow f$ uniformly on A . Then

$$\begin{aligned} &\mathbb{E}[\langle w(f_n(X)), Y - f_n(X) \rangle] \\ &= \mathbb{E}[\langle w(f_n(X)), Y - f_n(X) \rangle \mathbf{1}\{X \in A\}] + \mathbb{E}[\langle w(f_n(X)), Y - f_n(X) \rangle \mathbf{1}\{X \in A^c\}] \\ &\geq \mathbb{E}[\langle w(f_n(X)), Y - f_n(X) \rangle \mathbf{1}\{X \in A\}] - \mathbb{E}[\|Y - f_n(X)\| \mathbf{1}\{X \in A^c\}] \end{aligned} \quad (15.6.1)$$

because $\|w(s)\|_* \leq 1$. As $\|y - f_n(x)\| \mathbf{1}\{x \in A^c\} - \|y - f(x)\| \mathbf{1}\{x \in A^c\} \leq \|f(x) - f_n(x)\|$ by the triangle inequality, the last term in inequality (15.6.1) converges to $\mathbb{E}[\|Y - f(X)\| \mathbf{1}\{X \in A^c\}]$ as $n \rightarrow \infty$. Focusing on the first term in (15.6.1), for any $\epsilon_1 > 0$ the uniform convergence of f_n to f on A guarantees that for large enough n , we have

$$\begin{aligned} &\mathbb{E}[\langle w(f_n(X)), Y - f_n(X) \rangle] \\ &= \mathbb{E}[\langle w(f(X)), Y - f_n(X) \rangle \mathbf{1}\{X \in A\}] + \mathbb{E}[\langle w(f_n(X)) - w(f(X)), Y - f_n(X) \rangle \mathbf{1}\{X \in A\}] \\ &\geq \mathbb{E}[\langle w(f(X)), Y - f_n(X) \rangle \mathbf{1}\{X \in A\}] - k \sup_{x \in A} \|f(x) - f_n(x)\|_* \mathbb{E}[\|Y - f_n(X)\|] \\ &\geq \mathbb{E}[\langle w(f(X)), Y - f_n(X) \rangle \mathbf{1}\{X \in A\}] - \epsilon_1 \end{aligned}$$

Adding and subtracting $f(X)$ in the final expectation, we have

$$\begin{aligned}
& \mathbb{E}[\langle w(f(X)), Y - f_n(X) \rangle \mathbf{1}\{X \in A\}] \\
&= \mathbb{E}[\langle w(f(X)), Y - f(X) \rangle \mathbf{1}\{X \in A\}] + \mathbb{E}[\langle w(f(X)), f(X) - f_n(X) \rangle \mathbf{1}\{X \in A\}] \\
&\geq \mathbb{E}[\langle w(f(X)), Y - f(X) \rangle] - \mathbb{E}[\|Y - f(X)\| \mathbf{1}\{X \in A^c\}] - \mathbb{E}[\|f(X) - f_n(X)\|] \\
&\rightarrow \mathbb{E}[\langle w(f(X)), Y - f(X) \rangle] - \mathbb{E}[\|Y - f(X)\| \mathbf{1}\{X \in A^c\}].
\end{aligned}$$

Substituting these bounds into inequality (15.6.1), we have for any $\epsilon > 0$ that there exists a set A_ϵ with $P(A_\epsilon) \geq 1 - \epsilon$ and for which

$$\begin{aligned}
& \liminf_n \mathbb{E}[\langle w(f_n(X)), Y - f_n(X) \rangle] \\
&\geq \mathbb{E}[\langle w(f(X)), Y - f(X) \rangle] - 2\mathbb{E}[\|Y - f(X)\| \mathbf{1}\{X \in A_\epsilon^c\}].
\end{aligned}$$

For each $m \in \mathbb{N}$, let $B_m = \bigcup_{n \leq m} A_{1/n}$. Certainly $P(B_m) \geq 1 - 1/m$, and $f_n \rightarrow f$ uniformly on B_m (as the guarantees on $A_{1/n}$ from Egorov's theorem apply); the same argument thus gives

$$\begin{aligned}
& \liminf_n \mathbb{E}[\langle w(f_n(X)), Y - f_n(X) \rangle] \\
&\geq \mathbb{E}[\langle w(f(X)), Y - f(X) \rangle] - 2\mathbb{E}[\|Y - f(X)\| \mathbf{1}\{X \in B_m^c\}].
\end{aligned}$$

Because B_m is an increasing sequence of sets with $P(B_m) \geq 1 - 1/m$, the limit $B_\infty = \bigcup_m B_m$ satisfies $P(B_\infty) = 1$. For any $x \in B_\infty$, we see that $x \in B_m$ for some finite m ; trivially, for $x \in B_\infty$ we thus have $\|y - f(x)\| \mathbf{1}\{x \notin B_m\} \rightarrow \|y - f(x)\| \mathbf{1}\{x \notin B_\infty\} = 0$ as $m \rightarrow \infty$. Said differently, except on a null set, we have $\|y - f(x)\| \mathbf{1}\{x \notin B_m\} \rightarrow 0$ for P -almost all (x, y) , and this is certainly dominated by $\|y - f(x)\|$. Lebesgue's dominated convergence theorem then implies $\mathbb{E}[\|Y - f(X)\| \mathbf{1}\{X \notin B_m\}] \rightarrow 0$ as $m \rightarrow \infty$. Summarizing, we have shown that for any $w \in \mathcal{W}_k$, we have

$$\liminf_n \mathbb{E}[\langle w(f_n(X)), Y - f_n(X) \rangle] \geq \mathbb{E}[\langle w(f(X)), Y - f(X) \rangle].$$

By taking a supremum over $w \in \mathcal{W}_k$ in the last display and recognizing that $\epsilon > 0$ was arbitrary, we have shown that

$$\liminf_n \text{CE}(f_n, \mathcal{W}_k) \geq \text{CE}(f, \mathcal{W}_k)$$

for all $k < \infty$. By Lemma 15.6.3, for any integrable f and for each $\epsilon > 0$ there exists k such that

$$\sup_{w \in \mathcal{W}_k} \mathbb{E}[\langle w(f(X)), Y - f(X) \rangle] \geq \sup_{w \in \mathcal{W}} \mathbb{E}[\langle w(f(X)), Y - f(X) \rangle] - \epsilon.$$

and for this k we have

$$\liminf_n \text{CE}(f_n, \mathcal{W}_k) \geq \text{CE}(f, \mathcal{W}_k) \geq \text{CE}(f, \mathcal{W}) - \epsilon.$$

Noting that $\text{CE}(f_n, \mathcal{W}) \geq \text{CE}(f_n, \mathcal{W}_k)$ for any k and taking $\epsilon \rightarrow 0$ gives the lemma.

15.6.2 Proof of Proposition 15.5.2

The proof that $\text{CE}(\cdot, \mathcal{W}_{\|\cdot\|})$ identifies calibration (Definition 15.1, part (i)) is identical to the argument for Proposition 15.5.3, so we omit it.

Let $\mathcal{W} = \mathcal{W}_{\|\cdot\|}$ for shorthand, and consider a sequence of functions $f_n \rightarrow f$. Then

$$\text{CE}(f, \mathcal{W}) - \text{CE}(f_n, \mathcal{W}) \leq \sup_{w \in \mathcal{W}} \mathbb{E}[\langle w(f(X)), Y - f(X) \rangle - \langle w(f_n(X)), Y - f_n(X) \rangle]$$

and

$$\text{CE}(f_n, \mathcal{W}) - \text{CE}(f, \mathcal{W}) \leq \sup_{w \in \mathcal{W}} \mathbb{E}[\langle w(f_n(X)), Y - f_n(X) \rangle - \langle w(f(X)), Y - f(X) \rangle].$$

We focus on bounding the first display, as showing that the second tends to zero requires, *mutatis mutandis*, an identical argument.

Fix any $w \in \mathcal{W}$. Then

$$\begin{aligned} & \mathbb{E}[\langle w(f(X)), Y - f(X) \rangle - \langle w(f_n(X)), Y - f_n(X) \rangle] \\ &= \mathbb{E}[\langle w(f(X)) - w(f_n(X)), Y - f(X) \rangle] + \mathbb{E}[\langle w(f_n(X)), f_n(X) - f(X) \rangle] \\ &\leq \mathbb{E}[\min\{2, \|f(X) - f_n(X)\|\} \|Y - f(X)\|] + \mathbb{E}[\|f_n(X) - f(X)\|], \end{aligned}$$

where the inequality follows because $\|w(s) - w(s')\|_* \leq 2$ and $\|w(s) - w(s')\|_* \leq \|s - s'\|$ for any s, s' by construction. The second expectation certainly tends to zero as $n \rightarrow \infty$, so we consider the first. Define $g_n(x, y) = \min\{2, \|f(x) - f_n(x)\|\} \|y - f(x)\|$. Then $g_n(x, y) \leq g(x, y) = \|y - f(x)\|$, which has finite expectation by assumption. Moreover, Egorov's theorem (Lemma 15.6.1) guarantees that for each k , there is a set A_k with $P(A_k) \geq 1 - 1/k$ and for which $g_n \rightarrow 0$ uniformly on A_k (because $\mathbb{E}[\|f(X) - f_n(X)\|] \rightarrow 0$). Define $A_\infty = \bigcup_k A_k$, so that $P(A_\infty) = 1$, and $g_n(x, y) \rightarrow 0$ pointwise on A_∞ . Then the dominated convergence theorem guarantees that

$$\mathbb{E}[g_n(X, Y)] = \mathbb{E}[g_n(X, Y) \mathbf{1}\{(X, Y) \in A_\infty\}] + \underbrace{\mathbb{E}[g_n(X, Y) \mathbf{1}\{(X, Y) \notin A_\infty\}]}_{=0} \rightarrow 0.$$

Notably, this convergence is independent of w , and so we obtain

$$\limsup_n \{\text{CE}(f, \mathcal{W}) - \text{CE}(f_n, \mathcal{W})\} \leq 0.$$

A similar argument gives the converse bound.

15.6.3 Proof of Lemma 15.5.4

Define $f(s) = g(s) / \max\{1, \|g(s)\|_2\}$, so that $\mathbb{E}[\|g(s)\|_2] = \mathbb{E}[\langle f(s), g(s) \rangle]$. Using Lemma 15.6.3, we see that for each $n \in \mathbb{N}$ there exists a $C = C_n$ -Lipschitz function (where $C < \infty$) w_n with $\mathbb{E}[\|w_n(S) - f(S)\|] \leq \frac{1}{n}$, and w.l.o.g. we may assume $\|w_n(s)\|_2 \leq 1$ (by projection if necessary, which is Lipschitzian). Then

$$\mathbb{E}[\|g(S)\|_2] = \mathbb{E}[\langle f(S), g(S) \rangle] = \mathbb{E}[\langle f(S) - w_n(S), g(S) \rangle] + \underbrace{\mathbb{E}[\langle w_n(S), g(S) \rangle]}_{=0}.$$

Note that $w_n \rightarrow f$ in $L^1(P)$. Then for any $\epsilon > 0$, an application of Egorov's theorem (Lemma 15.6.1) and that $\mathbb{E}[\|g(S)\|] < \infty$ gives that we can find sets A_ϵ with $P(A_\epsilon) \geq 1 - \epsilon$ and for which $w_n \rightarrow f$ uniformly on A_ϵ . Then

$$\begin{aligned} \mathbb{E}[\|g(S)\|_2] &= \mathbb{E}[\langle f(S) - w_n(S), g(S) \rangle \mathbf{1}\{S \in A_\epsilon\}] + \mathbb{E}[\|g(S)\|_2 \mathbf{1}\{S \notin A_\epsilon\}] \\ &\leq \mathbb{E}\left[\sup_{s \in A_\epsilon} \|f(s) - w_n(s)\|_2 \|g(S)\|_2 \mathbf{1}\{S \in A_\epsilon\}\right] + \mathbb{E}[\|g(S)\|_2 \mathbf{1}\{S \notin A_\epsilon\}] \\ &\rightarrow \mathbb{E}[\|g(S)\|_2 \mathbf{1}\{S \notin A_\epsilon\}]. \end{aligned}$$

as $n \uparrow \infty$. We now employ the same device we use in the proof of Lemma 15.2.1. For $m \in \mathbb{N}$, let $B_m = \bigcup_{n \leq m} A_{1/n}$. Then $w_n \rightarrow f$ uniformly on B_m , and so $\mathbb{E}[\|g(S)\|_2] \leq \mathbb{E}[\|g(S)\|_2 \mathbf{1}\{S \notin B_m\}]$, that is, $\mathbb{E}[\|g(S)\|_2 \mathbf{1}\{S \in B_m\}] = 0$. Monotone convergence implies $0 = \lim_{m \rightarrow \infty} \mathbb{E}[\|g(S)\|_2 \mathbf{1}\{S \in B_m\}] = \mathbb{E}[\|g(S)\|_2 \mathbf{1}\{S \in B_\infty\}]$ where $B_\infty = \bigcup_m B_m$. As $P(B_\infty) = 1$ by continuity of measure, we have $\mathbb{E}[\|g(S)\|_2] = 0$, giving the lemma.

15.6.4 Proof of Theorem 15.5.6

The following lemma gives the lower bound in the theorem and is fairly straightforward.

Lemma 15.6.4. *For $S = f(X)$, we have*

$$p_{\text{cal},\text{up}}(f) \leq d_{\text{cal},\text{up}}(f) \leq \text{pce}(S). \quad (15.6.2)$$

Proof Fix any partition \mathcal{A} , and define $\mathbf{q}_{\mathcal{A}}(s)$ to be the (unique) set A such that $s \in A$ (so we quantize s). Then set $g(s) = \mathbb{E}[Y \mid S \in \mathbf{q}_{\mathcal{A}}(s)]$ to be the expectation of Y conditional on S being in the same partition element as s . Then $g(S) = \mathbb{E}[Y \mid g(S)]$ with probability 1, so that g is perfectly calibrated, and

$$\begin{aligned} p_{\text{cal},\text{up}}(f) &\leq d_{\text{cal},\text{up}}(f) \leq \mathbb{E}[\|S - g(S)\|] \\ &= \sum_{A \in \mathcal{A}} \mathbb{E}[\|S - \mathbb{E}[Y \mid S \in A]\| \mathbf{1}\{S \in A\}] \\ &\leq \sum_{A \in \mathcal{A}} \mathbb{E}[(\|S - \mathbb{E}[S \mid S \in A]\| + \|\mathbb{E}[S - Y \mid S \in A]\|) \mathbf{1}\{S \in A\}] \\ &\leq \sum_{A \in \mathcal{A}} \text{diam}(A) \mathbb{P}(S \in A) + \sum_{A \in \mathcal{A}} \|\mathbb{E}[(S - Y) \mathbf{1}\{S \in A\}]\|. \end{aligned}$$

Taking an infimum gives the claim (15.6.2). \square

To prove the claimed upper bound requires more work. For pedagogical reasons, let us attempt to prove a similar upper bound relating $\text{pce}(S)$ to $p_{\text{cal},\text{low}}(f)$. We might begin with a partition \mathcal{A} with maximal diameter $\text{diam}(A) \leq \epsilon$ for $A \in \mathcal{A}$, and for random variables (S, V, Y) , begin with the first term in the partition error, whence

$$\begin{aligned} \text{CE}(S, \mathcal{A}) &\leq \sum_{A \in \mathcal{A}} \|\mathbb{E}[(S - V) \mathbf{1}\{S \in A\}]\| + \|\mathbb{E}[(V - Y) \mathbf{1}\{S \in A\}]\| \\ &\leq \mathbb{E}[\|S - V\|] + \sum_{A \in \mathcal{A}} \|\mathbb{E}[(V - Y) \mathbf{1}\{V \in A\}]\| + \sum_{A \in \mathcal{A}} \|\mathbb{E}[(V - Y)(\mathbf{1}\{S \in A\} - \mathbf{1}\{V \in A\})]\| \\ &\leq \mathbb{E}[\|S - V\|] + \mathbb{E}[\|\mathbb{E}[Y \mid V] - V\|] + \sum_{A \in \mathcal{A}} \|\mathbb{E}[(V - Y)(\mathbf{1}\{S \in A\} - \mathbf{1}\{V \in A\})]\| \end{aligned}$$

by Jensen's inequality applied to conditional expectations, once we recognize $\mathbb{E}[(Y - V) \mathbf{1}\{V \in A\}] = \mathbb{E}[(\mathbb{E}[Y \mid V] - V) \mathbf{1}\{V \in A\}]$. For the final term, a straightforward computation yields

$$\begin{aligned} \sum_{A \in \mathcal{A}} \|\mathbb{E}[(V - Y)(\mathbf{1}\{S \in A\} - \mathbf{1}\{V \in A\})]\| &\leq \text{diam}(\mathcal{Y}) \sum_{A \in \mathcal{A}} [\mathbb{P}(S \in A, V \notin A) + \mathbb{P}(S \notin A, V \in A)] \\ &= 2 \text{diam}(\mathcal{Y}) \mathbb{P}(S \text{ and } V \text{ belong to different } A \in \mathcal{A}). \end{aligned}$$

If S and V had continuous distributions, we would expect the probability that they fail to belong to the same partition elements to scale as $\mathbb{E}[\|S - V\|]$. This may fail, but to rectify the issue, we can randomize.

Consequently, let us consider the *randomized partition error*, which we index with $\varepsilon > 0$ and for $U \sim \text{Uniform}[-1, 1]^k$ define as

$$\text{rpce}_\varepsilon(S) := \inf_{\mathcal{A}} \left\{ \sum_{A \in \mathcal{A}} \|\mathbb{E}[(S - Y)\mathbf{1}\{S + \varepsilon U \in A\}]\| + \sum_{A \in \mathcal{A}} \text{diam}(A)\mathbb{P}(S \in A) \right\}. \quad (15.6.3)$$

(The choice of uniform $[-1, 1]^k$ is only made for convenience in the calculations to follow.) Letting $c_k = \|\mathbf{1}_k\|_*$, we see immediately that

$$\text{pce}(S) \leq \text{rpce}_\varepsilon(S) + c_k \varepsilon$$

for all $\varepsilon \geq 0$. We can say more.

Lemma 15.6.5. *Let $\varepsilon > 0$. Then for any random variable V ,*

$$\text{rpce}_\varepsilon(S) \leq \mathbb{E}[\|S - V\|] + \mathbb{E}[\|\mathbb{E}[Y | V] - V\|] + \frac{2k}{\varepsilon} \mathbb{E}[\|Y - S\| \|V - S\|_\infty].$$

Note that by combining Lemma 15.6.5 with the display above and recognizing that $\|Y - S\| \leq \text{diam}(\mathcal{Y})$ with probability 1, we have the theorem.

Proof We replicate the calculation bounding $\text{CE}(S, \mathcal{A})$ above, but while allowing the randomization. Let \mathcal{A} be a partition of \mathbb{R}^k into hypercubes of width ε , that is, $[-\varepsilon, \varepsilon]^k + \varepsilon z$, where $z \in 2\mathbb{Z}^k$ ranges over integer vectors with even entries. Then $\text{diam}(A) \leq c_k \varepsilon$, and

$$\begin{aligned} & \|\mathbb{E}[(S - Y)\mathbf{1}\{S + \varepsilon U \in A\}]\| \\ & \leq \|\mathbb{E}[(S - V)\mathbf{1}\{S + \varepsilon U \in A\}]\| + \|\mathbb{E}[(V - Y)\mathbf{1}\{V + \varepsilon U \in A\}]\| \\ & \quad + \|\mathbb{E}[(V - Y)(\mathbf{1}\{S + \varepsilon U \in A\} - \mathbf{1}\{V + \varepsilon U \in A\})]\| \\ & \leq \|\mathbb{E}[(S - V)\mathbf{1}\{S + \varepsilon U \in A\}]\| + \|\mathbb{E}[(V - Y)\mathbf{1}\{V + \varepsilon U \in A\}]\| \\ & \quad + \mathbb{E}[\|V - Y\| \cdot (\mathbb{P}(V + \varepsilon U \in A, S + \varepsilon U \notin A | V, S) + \mathbb{P}(S + \varepsilon U \in A, V + \varepsilon U \notin A | V, S, Y))] \end{aligned}$$

Summing over sets A and using the triangle inequality and that $S + \varepsilon U \in A$ for some A , we find

$$\begin{aligned} \sum_{A \in \mathcal{A}} \|\mathbb{E}[(S - Y)\mathbf{1}\{S + \varepsilon U \in A\}]\| & \leq \mathbb{E}[\|S - V\|] + \mathbb{E}[\|\mathbb{E}[Y | V] - V\|] \\ & \quad + 2\mathbb{E}\left[\|V - Y\| \sum_{A \in \mathcal{A}} \mathbb{P}(V + \varepsilon U \in A, S + \varepsilon U \notin A | V, S, Y)\right]. \end{aligned} \quad (15.6.4)$$

We now may bound the probability in inequality (15.6.4). Recall that $A = [-\varepsilon, \varepsilon]^k + \varepsilon z$ for some $z \in 2\mathbb{Z}^k$, and fix $v, s \in \mathbb{R}^k$. Let $\mathbb{B} = [-1, 1]^k$ be the ℓ_∞ ball. Then

$$\begin{aligned} \mathbb{P}(v + \varepsilon U \in \mathbb{B}, s + \varepsilon U \notin \mathbb{B}) & = \mathbb{P}(U \notin \varepsilon^{-1}(\mathbb{B} - s) | U \in \varepsilon^{-1}(\mathbb{B} - v))\mathbb{P}(v + \varepsilon U \in \mathbb{B}) \\ & \leq \frac{k \|s - v\|_\infty}{\varepsilon} \mathbb{P}(v + \varepsilon U \in \mathbb{B}), \end{aligned} \quad (15.6.5)$$

where inequality (15.6.5) follows because if $s, v \in \mathbb{R}^k$ are the centers of two ℓ_∞ balls B_s and B_v of radius 1, and if $\delta = \|s - v\|_\infty$, then the volume of $B_v \setminus B_s$ is at most $k\delta^k/\delta^{k-1} = k\delta$. (See

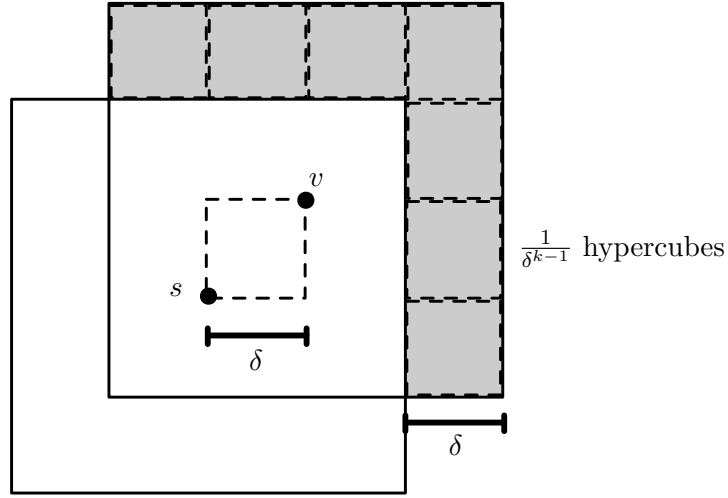


Figure 15.2. The volume argument in inequality (15.6.5). In k dimensions, the hypercube of side-length δ can be replicated $1/\delta^{k-1}$ times on each exposed base of the cube centered at v , where $\delta = \|s - v\|_\infty$. There are at most k such faces, giving volume at most $k\delta^k/\delta^{k-1} = k\delta$ to the gray region.

Figure 15.2. The k -dimensional surface area of one side of a hypercube of radius δ is $2k\delta^{k-1}$, and we can put at most $1/\delta^{k-1}$ boxes in each facial part of the grey region.)

Substituting inequality (15.6.5) into the bound (15.6.4) and conditioning and deconditioning on V, S , we find that

$$\begin{aligned}
 & \sum_{A \in \mathcal{A}} \|\mathbb{E}[(S - Y)\mathbf{1}\{S + \varepsilon U \in A\}]\| \\
 & \leq \mathbb{E}[\|S - V\|] + \mathbb{E}[\|\mathbb{E}[Y | V] - V\|] + \frac{2k}{\varepsilon} \mathbb{E} \left[\|V - Y\| \sum_{A \in \mathcal{A}} \|V - S\|_\infty \mathbf{1}\{V + \varepsilon U \in A\} \right] \\
 & = \mathbb{E}[\|S - V\|] + \mathbb{E}[\|\mathbb{E}[Y | V] - V\|] + \frac{2k}{\varepsilon} \mathbb{E}[\|Y - V\| \|V - S\|_\infty].
 \end{aligned}$$

Taking an infimum over partitions \mathcal{A} gives the lemma. \square

15.7 Bibliography

Draft: Calibration remains an active research area. The initial references for online calibration are Foster and Vohra [90], Dawid and Vovk [64]. The idea of calibrating is most present in Foster and Hart [94]. Our proof of calibrating is based on Kumar et al. [128]. Blasiok et al. [34] demonstrate the equivalence of the different metrics for measuring calibration, focusing on the case of binary prediction; the extension to vector-valued Y appears to be new. The ideas of the postprocessing gap and also descend from Blasiok et al. [35], and the connections with general proper losses also appear to be new. Propositions 15.5.1, 15.5.2, and 15.5.3 are new in that they are the first to demonstrate that the measures are valid calibration measures (Definition 15.1, part (i)).

JCD Comment: A few more things to add either in the bibliography or the introduction to the section:

1. We only really do calibration for binary/multiclass things. One would also really like to predict full distributions P_t on general outcomes Y , which is harder (nearly impossible) to do in any conditional sense.
2. It's much easier to do predictive inference (cover) because don't need accuracy
3. Maybe comment on variants for top entry (from multiclass to binary) classification and why that is important. Maybe in the middle, maybe here.

15.8 Exercises

Exercise 15.1: We say \mathcal{W} is *weakly complete* if $\mathbb{E}[\langle w(S), g(S) \rangle] = 0$ implies that $g(S) = 0$ with probability 1. Let \mathcal{W} be any symmetric weakly complete collection of functions. Show that

$$\text{CE}(f, \mathcal{W})$$

is a sound and complete (15.2.2) calibration measure.

Exercise 15.2 (Marginal calibration): A predictor $f : \mathcal{X} \rightarrow \mathbb{R}^k$ of $Y \in \mathbb{R}^k$ is *per-class calibrated* if for each $j \in \{1, \dots, k\}$ we have $\mathbb{E}[Y_j | f_j(X)] = f_j(X)$ with probability 1.

- (a) Show that if $f : \mathcal{X} \rightarrow \mathbb{R}^k$ is calibrated, then it is marginally calibrated.
- (b) Let \mathcal{W} be symmetric and weakly complete (Exercise 15.1) and $\mathcal{W}^k = \{(w_1, \dots, w_k) \mid w_j \in \mathcal{W}\}$ be the vector-valued functions with components in \mathcal{W} , where $w \in \mathcal{W}^k$ satisfies $w(s) = (w_1(s), \dots, w_k(s))$ for $s \in \mathbb{R}^k$. Show that $\text{CE}(f, \mathcal{W}^k)$ is sound and complete for marginal calibration, that is, $\text{CE}(f, \mathcal{W}^k) = 0$ if and only if f is marginally calibrated.

Exercise 15.3 (Top-K calibration): In k -class multiclass classification, where we treat $y \in \{e_1, \dots, e_k\} \subset \mathbb{R}^k$, instead of obtaining calibration of an entire vector $\mathbb{E}[Y | f(X)]$ we can ask for calibration of the top K predictions. Let $f_{(1)}(x) \geq f_{(2)}(x) \geq \dots \geq f_{(k)}(x)$ denote the sorted elements of $f(x) \in \mathbb{R}^k$. Then if $f_j(X)$ is among the K -highest predictions, that is, $f_j(X) \geq f_{(K)}(X)$, we ask that it be calibrated,

$$\mathbb{P}(Y = e_j \mid f_j(X), f_j(X) \geq f_{(K)}(X)) = f_j(X).$$

For a function class \mathcal{W} mapping $\mathbb{R} \rightarrow \mathbb{R}$, define the function class

$$\mathcal{W}_{\text{top-}K} := \left\{ w(s) = \sum_{j=1}^k w_j(s_j) e_j \mathbf{1}\{s_j \geq s_{(K)}\} \mid w_j \in \mathcal{W}, j = 1, \dots, k \right\}.$$

Show that if \mathcal{W} is weakly complete (Exercise 15.1) and symmetric, then $\text{CE}(f, \mathcal{W}_{\text{top-}K})$ is sound and complete for top- K calibration.

JCD Comment: Add a uniform convexity version of Proposition 15.3.5 as an exercise.

JCD Comment: Can we add an exercise about achieving weak calibration for different classes of functions?

JCD Comment: A few potential exercises:

- (i) Deal with any class \mathcal{W} for which $\mathbb{E}[\langle w, f \rangle] = 0$ for all $w \in \mathcal{W}$ means $f = 0$, then still get a continuous calibration measure

JCD Comment: Exercise: do Aaditya's top-class calibration approach.

JCD Comment: Do we need more commentary on calibrating? Maybe an exercise on empirics? Project ideas: calibrating with witnesses in higher dimensions, doing calibrating in higher dimensions, optimality results / lower bounds.

JCD Comment: Do Example 3.2 of Kumar et al. [128] as exercise

JCD Comment: Coding and empirical exercises on calibration?

JCD Comment: Remark on impossibility of inference of ece? Exercises on its impossibility too, perhaps, and one-sided estimation of it. And maybe some minimax lower bounds on the Lipschitz one as well I think.

JCD Comment: Exercise potential: let \mathcal{W} be a collection from an RKHS

Chapter 16

Classification, Divergences, and Surrogate Risk

While proper losses—a major focus in chapters 14 and 15—provide a principled approach to choosing a loss and prediction the “right” distribution in prediction problems, they are not a panacea. In many contexts, we wish to focus more exclusively on optimal prediction of a particular target Y rather than its distribution, whether because the problem at hand does not require distributional predictions or because Y itself is so complex that a distribution over Y would be generally intractable to describe. For example, in *structured prediction* problems, the target Y may be a complex or combinatorial object, such as labeling each pixel in an image as foreground or background, a parse tree of a sentence in natural language processing, or a ranking of items.

The setting for much of this chapter has similarities to that in the preceding two: we assume data coming in pairs $(X, Y) \in \mathcal{X} \times \mathcal{Y}$, but now we abstractly consider a function $f : \mathcal{X} \rightarrow \mathbb{R}^k$ for some k , and we seek the function f minimizing the risk, or expected loss,

$$R(f) := \mathbb{E}[\ell(f(X), Y)],$$

where $\ell : \mathbb{R}^k \times \mathcal{Y} \rightarrow \mathbb{R}$ is a loss function. Typically, we think of some decoder or prediction function that transforms f into a prediction $\hat{y}(f(x))$ of y . Binary and multiclass classification provide the cleanest motivation of this approach. In the former, we consider $Y \in \{-1, 1\}$, take $f : \mathcal{X} \rightarrow \mathbb{R}$ with prediction $\hat{y} = \text{sign}(f(x))$, and the natural loss is the zero-one error

$$\ell_{0-1}(s, y) = \mathbf{1}\{sy \leq 0\}$$

for $s \in \mathbb{R}$. In the latter, for multiclass classification, we take $Y \in \{1, \dots, k\}$ and $f : \mathcal{X} \rightarrow \mathbb{R}^k$ with prediction $\hat{y} = \arg\max_j f_j(x)$, and the zero-one error for $s \in \mathbb{R}^k$ and label $y \in [k]$ becomes

$$\ell_{0-1}(s, y) = \mathbf{1}\left\{s_y \leq \max_{j \neq y} s_j\right\},$$

which is 0 when the score s_y assigned to the correct label y is greater than all others.

In many of the cases we consider in this chapter, the loss ℓ is hard to minimize: it is non-convex and, as in the case of the zero-one losses above, non-differentiable, NP-hard to minimize in general, and gradient information does not permit even some reasonable local search for a predictive function. As a consequence, this chapter takes a different tack, where we search for

surrogate losses $\varphi : \mathbb{R}^k \times \mathcal{Y} \rightarrow \mathbb{R}$ that admit nicer properties—such as convexity and smoothness—such that minimizing φ is equivalent in some sense to minimizing the original loss ℓ . To define our desiderata, let

$$R_\varphi(f) := \mathbb{E}[\varphi(f(X), Y)]$$

be the surrogate risk of a function f , and define the *Bayes risk* (for ℓ and φ) as the minimal possible expected loss,

$$R^* := \inf_f \mathbb{E}[\ell(f(X), Y)] \quad \text{and} \quad R_\varphi^* := \inf_f \mathbb{E}[\varphi(f(X), Y)],$$

where the infima are over the class of all measurable functions. We then seek *surrogate risk consistency*, essentially, a type of infinite sample consistency, sometimes called *Fisher consistency*, a population-level guarantee that minimizing R_φ guarantees minimizing R .

Definition 16.1 (Surrogate risk consistency). *The loss φ is surrogate risk consistent for the loss ℓ if for any sequence of functions f_n and any distribution on (X, Y) ,*

$$R_\varphi(f_n) \rightarrow R_\varphi^* \quad \text{implies} \quad R(f_n) \rightarrow R^*.$$

This chapter develops the theory of such surrogate risk consistency. The theory obtains its cleanest and most transparent form in the case of binary classification with the zero-one error, but it extends beyond this, including to multiclass and structured prediction problems. As we will show, the dualities we develop in Chapter 14 and connections with generalized entropies remain important; in many cases, any loss φ generating a suitably concave generalized entropy guarantees consistency. Deeper equivalence results between loss functions arise via this entropy and information-based perspective, and we close the chapter by developing these. Throughout, we will elide measure-theoretic details, as they are secondary to the main thrusts of the arguments (but see the bibliographic section).

16.1 Surrogate risk consistency in binary classification

The first step, both because it offers the most complete solution and because the ideas are most transparent, is to consider consistency for binary classification. In this case, we take $Y \in \{-1, 1\}$, and consider margin-based losses, where the goal is to find a predictor $f : \mathcal{X} \rightarrow \mathbb{R}$ predicting the correct sign of y . Thus, for $s \in \mathbb{R}$, we have the zero-one loss

$$\ell(s, y) = \mathbf{1}\{sy \leq 0\} \quad \text{and} \quad \varphi(s, y) = \phi(y s),$$

where $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is a convex function acting on the *margin* of a predictor, where the margin of f on an example (x, y) is

$$yf(x).$$

A large margin $yf(x) \gg 0$ thus corresponds to a correct prediction, while non-positive margin $yf(x) \leq 0$ means that $f(x)$ has opposite sign to y .

Example 16.1.1 (Common convex loss functions): Several convex losses are frequent in the literature on classification. These include the logistic loss,

$$\phi(t) = \log(1 + e^{-t}),$$

the exponential loss, where

$$\phi(t) = e^{-t},$$

the squared error

$$\phi(t) = (1 - t)^2,$$

and the hinge loss and squared hinge loss

$$\phi(t) = [1 - t]_+ \quad \text{and} \quad \phi(t) = [1 - t]_+^2.$$

Each of these is nonnegative and satisfies $\phi'(0) < 0$. \diamond

The key step to understanding surrogate-risk consistency is to move from the full population expectation to conditional expectations given X : as we work with arbitrary measurable functions $f : \mathcal{X} \rightarrow \mathbb{R}$, we can essentially choose $f(X)$ to be optimally predict of Y given X . To that end, let

$$\eta(x) := \mathbb{P}(Y = 1 \mid X = x)$$

be the conditional probability that $Y = 1$ given $X = x$, so that

$$R(f) = \mathbb{P}(\text{sign}(f(X)) \neq Y) = \mathbb{E}[\eta(X)\mathbf{1}\{f(X) \leq 0\} + (1 - \eta(X))\mathbf{1}\{f(X) \geq 0\}]$$

and

$$R_\varphi(f) = \mathbb{E}[\eta(X)\phi(f(X)) + (1 - \eta(X))\phi(-f(X))].$$

Immediately, we see that $f^*(x) = \eta(x) - \frac{1}{2}$ minimizes R (one may arbitrarily modify f^* on the set $\{x \mid \eta(x) = \frac{1}{2}\}$). Defining the conditional risks

$$r(s, \eta) = \eta\mathbf{1}\{s \leq 0\} + (1 - \eta)\mathbf{1}\{s \geq 0\} \quad \text{and} \quad r_\phi(s, \eta) = \eta\phi(s) + (1 - \eta)\phi(-s)$$

we evidently have

$$R(f) = \mathbb{E}[r(f(X), \eta(X))] \quad \text{and} \quad R_\varphi(f) = \mathbb{E}[r_\phi(f(X), \eta(X))].$$

Thus, we seek relationships that relate $r(s, \eta)$ and $r_\phi(s, \eta)$ to one another that guarantee consistency.

Example 16.1.2 (Exponential loss): Consider the exponential loss, which AdaBoost and other boosting algorithms use, which sets $\phi(s) = e^{-s}$. In this case, for $\eta \in (0, 1)$ we have

$$\operatorname{argmin}_s r_\phi(s, \eta) = \frac{1}{2} \log \frac{\eta}{1 - \eta} \quad \text{because} \quad \frac{\partial}{\partial s} r_\phi(s, \eta) = -\eta e^{-s} + (1 - \eta)e^s.$$

Solving for s by setting the derivative to zero yields $e^{2s} = \frac{\eta}{1 - \eta}$, so that $f_\phi^*(x) = \frac{1}{2} \log \frac{\eta(x)}{1 - \eta(x)}$ minimizes $R_\varphi(f)$, where we allow f^* to take infinite values if $\eta \in \{0, 1\}$. As $\text{sign}(f_\phi^*(x)) = \text{sign}(2\eta(x) - 1)$, this is optimal for the zero-one loss as well. \diamond

To provide a more quantitative version of this pointwise identification of minimizers, we compare the conditional ϕ -risks. Intuitively, if the s minimizing $r_\phi(s, \eta)$ has the same sign as $2\eta - 1$, then we should obtain consistency. Said differently, if one cannot minimize r_ϕ accurately when we constrain s to have the incorrect sign, then any minimizer of R_φ should agree with the minimizers of R . Following this intuition, define the the infimal conditional ϕ -risks as

$$r_\phi^*(\eta) := \inf_s r_\phi(s, \eta) \quad \text{and} \quad r_\phi^{\text{wrong}}(\eta) := \inf_{s(\eta - 1/2) \leq 0} r_\phi(s, \eta).$$

Then we expect that $r_\phi^*(\eta) < r_\phi^{\text{wrong}}(\eta)$ for all $\eta \neq \frac{1}{2}$ is sufficient to guarantee surrogate risk consistency for the zero-one loss.

To make this intuition rigorous, define the sub-optimality function $\Delta_\phi : [0, 1] \rightarrow \mathbb{R}$

$$\Delta_\phi(\delta) := r_\phi^{\text{wrong}}\left(\frac{1+\delta}{2}\right) - r_\phi^*\left(\frac{1+\delta}{2}\right). \quad (16.1.1)$$

We may now define

Definition 16.2 (Classification calibration). *The margin-based loss ϕ is classification calibrated if $\Delta_\phi(\delta) > 0$ for all $\delta > 0$. Equivalently, for any $\eta \neq \frac{1}{2}$, we have $r_\phi^*(\eta) < r_\phi^{\text{wrong}}(\eta)$.*

Examples can help to visualize this definition.

Example (Example 16.1.2 continued): For the exponential loss, we have

$$r_\phi^{\text{wrong}}(\eta) = \inf_{s(2\eta-1) \leq 0} \{\eta e^{-s} + (1-\eta)e^s\} = e^0 = 1$$

while the unconstrained minimal conditional risk is

$$r_\phi^*(\eta) = \eta \sqrt{\frac{1-\eta}{\eta}} + (1-\eta) \sqrt{\frac{\eta}{1-\eta}} = 2\sqrt{\eta(1-\eta)},$$

so that $\Delta_\phi(\delta) = 1 - \sqrt{1-\delta^2} \geq \frac{1}{2}\delta^2$, where we have used that $\sqrt{a+x} \leq \sqrt{a} + \frac{x}{2\sqrt{a}}$ for all $a > 0, x \in \mathbb{R}$. \diamond

Example 16.1.3 (Hinge loss): We can also consider the hinge loss $\phi(s) = [1-s]_+$. Computing the minimizers of the conditional risk, we have

$$r_\phi(s, \eta) = \eta [1-s]_+ + (1-\eta) [1+s]_+,$$

whose unique minimizer (for $\eta \notin \{0, \frac{1}{2}, 1\}$) is $s(\eta) = \text{sign}(2\eta - 1)$, while $s(\eta) = \text{sign}(2\eta - 1)$ is a minimizer for all η . We thus have

$$r_\phi^*(\eta) = 2 \min\{\eta, 1-\eta\} \quad \text{and} \quad r_\phi^{\text{wrong}}(\eta) = \eta + (1-\eta) = 1.$$

We obtain $\Delta_\phi(\delta) = 1 - \min\{1+\delta, 1-\delta\} = \delta$. \diamond

Example 16.1.4 (Squared error): Let the margin-based loss be $\phi(s, y) = \frac{1}{2}(s-y)^2$. Then we have

$$r_\phi^*(\eta) = \inf_s \left\{ \frac{\eta}{2}(s-1)^2 + \frac{1-\eta}{2}(s+1)^2 \right\} = 2\eta(1-\eta)$$

because $s(\eta) = 2\eta - 1$ minimizes $r_\phi(s, \eta)$. On the other hand, $r_\phi^{\text{wrong}}(\eta) = \frac{1}{2}$ for all η , and so we have $\Delta_\phi(\delta) = \frac{1}{2} - \frac{1}{2}(1-\delta)(1+\delta) = \frac{1}{2}\delta^2$. \diamond

JCD Comment: Put figures here with exponential loss and conditional probability η .

16.1.1 A general classification calibration result

In each of the preceding examples, we have seen that the suboptimality function Δ_ϕ is convex. While this need not hold generally, by a careful application of Jensen's inequality we can leverage a convexified variant of Δ_ϕ yields a convex function ψ , strictly positive on \mathbb{R}_{++} , such that

$$\psi(R_\phi(f) - R_\phi^*) \leq R(f) - R^*,$$

where we have used $R_\phi(f) = \mathbb{E}[\phi(Yf(X))]$ to make clear that we use the margin-based loss $\varphi(s, y) = \phi(y s)$. The key is to leverage the biconjugate of Δ_ϕ , which (recall Lemma 14.1.1) is the largest (closed) convex function below Δ_ϕ . To that end, we define the ψ -transform

$$\psi_\phi(\delta) := \Delta_\phi^{**}(\delta) \quad (16.1.2)$$

of the margin-based loss ϕ . With this, we obtain the following characterization of classification calibration.

Theorem 16.1.5. *Let ϕ be a margin-based loss function and ψ the associated ψ -transform. Then for any $f : \mathcal{X} \rightarrow \mathbb{R}$,*

$$\psi(R(f) - R^*) \leq R_\phi(f) - R_\phi^*. \quad (16.1.3)$$

Moreover, if ϕ is nonnegative, the following three are equivalent:

(i) The loss ϕ is classification-calibrated (Definition 16.2).

(ii) For any sequence $\delta_n \in [0, 1]$,

$$\psi(\delta_n) \rightarrow 0 \quad \text{if and only if} \quad \delta_n \rightarrow 0.$$

(iii) For any sequence of measurable functions $f_n : \mathcal{X} \rightarrow \mathbb{R}$,

$$R_\phi(f_n) \rightarrow R_\phi^* \quad \text{implies} \quad R(f_n) \rightarrow R^*.$$

The proof of the theorem relies on the convex analysis we develop in the brief primer in Section 14.1.1 and Appendix B, so we defer it to Section 16.1.3.

Theorem 16.1.5 provides concrete conditions to guarantee infinite-sample consistency of a margin-based surrogate loss $\varphi(s, y) = \phi(y s)$. If the associated suboptimality gap $\Delta_\phi(\delta)$ is strictly positive for $\delta > 0$, then the associated ψ -transform (16.1.2) is also strictly positive, and we have the quantitative guarantee that

$$\psi_\phi(R(f) - R^*) \leq R_\phi(f) - R_\phi^*,$$

so that a small excess surrogate risk guarantees small excess zero-one error, though different loss functions yield different explicit upper bounds.

By recalling Examples 16.1.2, 16.1.3, and 16.1.4, we can see immediate applications of the theorem. For each of these, we have that $\psi(\delta) = \Delta(\delta)$, as the function Δ is convex. Considering each in turn, for the exponential loss $\phi(s) = \exp(-s)$, we have $\Delta(\delta) = 1 - \sqrt{1 - \delta^2} = \frac{1}{2}\delta^2 + O(\delta^4) \geq \frac{1}{2}\delta^2$, while for the hinge loss $\phi(s) = [1 - s]_+$, we have $\Delta(\delta) = \delta$. Thus we obtain for any f that

$$\mathbb{P}(Yf(X) \leq 0) - \inf_f \mathbb{P}(Yf(X) \leq 0) \leq \mathbb{E} [[1 - Yf(X)]_+] - \inf_f \mathbb{E} [[1 - Yf(X)]_+].$$

On the other hand, for the exponential loss, we have

$$\frac{1}{2} \left(\mathbb{P}(Yf(X) \leq 0) - \inf_f \mathbb{P}(Yf(X) \leq 0) \right)^2 \leq \mathbb{E}[\exp(-Yf(X))] - \inf_f \mathbb{E}[\exp(-Yf(X))],$$

so that the bound on the zero-one error that the exponential loss guarantees is quadratically worse than that the hinge loss provides. The squared loss $\frac{1}{2}(f(x) - y)^2$ has suboptimality function $\Delta(\delta) = \frac{1}{2}\delta^2$, yielding a similar guarantee as that the exponential loss provides while also guaranteeing that regressing directly on ± 1 labels is consistent.

16.1.2 Convex losses for binary classification

Most commonly, one takes a convex margin-based loss ϕ , because these admit computationally efficient minimization and fitting procedures. We can thus specialize Theorem 16.1.5 to the convex case, obtaining a much simpler characterization of classification calibrated losses. Intuitively, if $\phi'(0) < 0$, then any scalar minimizing the conditional ϕ -risk

$$r_\phi(s, \eta) = \eta\phi(s) + (1 - \eta)\phi(-s)$$

should necessarily satisfy $\text{sign}(s) = \text{sign}(2\eta - 1)$. This is precisely the correct condition.

JCD Comment: Put a figure, including the fact that $r'_\phi(0, \eta) = (2\eta - 1)\phi'(0) < 0$ when $\eta > \frac{1}{2}$, so minimizer must be to the right.

Theorem 16.1.6. *If ϕ is convex, then ϕ is classification calibrated if and only if $\phi'(0)$ exists and $\phi'(0) < 0$.*

Proof Let $\eta > \frac{1}{2}$ w.l.o.g. First, suppose that ϕ is differentiable at 0 and $\phi'(0) < 0$. Then the dom ϕ necessarily includes an interval around 0, and

$$r_\phi(s, \eta) = \eta\phi(s) + (1 - \eta)\phi(-s)$$

satisfies $r'_\phi(0, \eta) = (2\eta - 1)\phi'(0)$; if $\phi'(0) < 0$, this quantity is negative for $\eta > 1/2$. Thus the minimizing $s(\eta) \in (0, \infty]$. Indeed, we have

$$r_\phi(s, \eta) \geq r_\phi(0, \eta) + (2\eta - 1)\phi'(0)s$$

by the first-order inequality for convexity, and the final term is strictly positive for $s < 0$. Moreover,

$$r_\phi(s, \eta) = r_\phi(0, \eta) + (2\eta - 1)\phi'(0)s + o(s) < r_\phi(0, \eta) = \phi(0)$$

for $s > 0$ but near 0. Thus $r_\phi^{\text{wrong}}(\eta) = r_\phi(0, \eta) = \phi(0)$, while $r_\phi^*(\eta) < \phi(0)$ for all $\eta \neq \frac{1}{2}$. In particular, $\Delta_\phi(\delta) > 0$ for $\delta > 0$ as desired.

To see the converse, we must prove that ϕ is differentiable at 0. Recall that a subgradient g_s of the function ϕ at $s \in \mathbb{R}$ is any g_s such that $\phi(t) \geq \phi(s) + g_s(t - s)$ for all $t \in \mathbb{R}$. Because ϕ is classification calibrated, its domain necessarily includes an interval around 0; thus that ϕ is subdifferentiable at 0 and the subgradient set $\partial\phi(0) \neq \emptyset$ (See Theorem B.3.3 or Proposition B.3.20 in Appendix B.) Let $g_1, g_2 \in \partial\phi(0)$. We show that both $g_1, g_2 < 0$ and $g_1 = g_2$, implying that ϕ is differentiable at 0.

By convexity we have

$$\begin{aligned} r_\phi(s, \eta) &\geq \eta(\phi(0) + g_1 s) + (1 - \eta)(\phi(0) - g_2 s) \\ &= [\eta g_1 - (1 - \eta)g_2] s + \phi(0). \end{aligned} \quad (16.1.4)$$

We first show that $g_1 = g_2$, meaning that ϕ is differentiable. Without loss of generality, assume $g_1 > g_2$. Then for $\eta > 1/2$, we would have $\eta g_1 - (1 - \eta)g_2 > 0$, which would imply that

$$r_\phi(s, \eta) \geq \phi(0) \geq \inf_{s' \leq 0} \{\eta \phi(s') + (1 - \eta)\phi(-s')\} = r_\phi^{\text{wrong}}(\eta),$$

for all $s \geq 0$ by (16.1.4), by taking $s' = 0$ in the second inequality. By our assumption of classification calibration, for $\eta > 1/2$ we know that

$$\inf_s r_\phi(s, \eta) < \inf_{s \leq 0} r_\phi(s, \eta) = r_\phi^{\text{wrong}}(\eta) \quad \text{so} \quad r_\phi^*(\eta) = \inf_{s \geq 0} r_\phi(s, \eta),$$

and under the assumption that $g_1 > g_2$ we obtain $r_\phi^*(\eta) = \inf_{s \geq 0} r_\phi(s, \eta) > r_\phi^{\text{wrong}}(\eta)$, which is a contradiction to classification calibration. We thus obtain $g_1 = g_2$, so that the function ϕ has a unique subderivative at $s = 0$ and is thus differentiable.

Now that we know ϕ is differentiable at 0, consider

$$\eta \phi(s) + (1 - \eta)\phi(-s) \geq (2\eta - 1)\phi'(0)s + \phi(0).$$

If $\phi'(0) \geq 0$, then for $s \geq 0$ and $\eta > 1/2$ the right hand side must be at least $\phi(0)$, which contradicts classification calibration, because $r_\phi^*(\eta) < r_\phi^{\text{wrong}}(\eta)$ as in the preceding argument. \square

Theorem 16.1.6 makes it easy to determine whether a convex margin-based loss is classification calibrated: simply take its derivative at 0. Each of Examples 16.1.2–16.1.4 is thus immediately classification calibrated. Other examples follow immediately as well.

Example 16.1.7 (Logistic loss): The logistic loss $\phi(t) = \log(1 + e^{-t})$ satisfies $\phi'(0) = -\frac{1}{2}$ and so is classification calibrated. \diamond

Example 16.1.8 (Squared hinge loss): The squared hinge loss $\phi(t) = [1 - t]_+^2$ satisfies $\phi'(0) = -2$ and so is classification calibrated. \diamond

16.1.3 Proof of Theorem 16.1.5

Throughout the proof, as the margin-based loss ϕ remains consistent, we let $\Delta = \Delta_\phi$ and $\psi = \psi_\phi$. We begin with the first statement (16.1.3). By a calculation, the gap (for a fixed margin s) in conditional 0-1 risk is

$$\begin{aligned} r(s, \eta) - \inf_s r(s, \eta) &= \eta \mathbf{1}\{s \leq 0\} + (1 - \eta) \mathbf{1}\{s \geq 0\} - \eta \mathbf{1}\{\eta \leq 1/2\} - (1 - \eta) \mathbf{1}\{\eta \geq 1/2\} \\ &= \begin{cases} 0 & \text{if } \text{sign}(s) = \text{sign}(\eta - \frac{1}{2}) \\ \eta \vee (1 - \eta) - \eta \wedge (1 - \eta) = |2\eta - 1| & \text{if } \text{sign}(s) \neq \text{sign}(\eta - \frac{1}{2}). \end{cases} \end{aligned}$$

In particular, we obtain the gap in risks

$$R(f) - R^* = \mathbb{E}[\mathbf{1}\{\text{sign}(f(X)) \neq \text{sign}(2\eta(X) - 1)\} |2\eta(X) - 1|]. \quad (16.1.5)$$

We use expression (16.1.5) to upper bound $R(f) - R^*$ via the ϕ -risk, for which we use the ψ -transform (16.1.2). By Jensen's inequality,

$$\psi(R(f) - R^*) \leq \mathbb{E}[\psi(\mathbf{1}\{\text{sign}(f(X)) \neq \text{sign}(2\eta(X) - 1)\} |2\eta(X) - 1|)].$$

We always have $\Delta(0) = 0$ as $r_\phi^{\text{wrong}}(1/2) = \inf_{s(1-1) \leq 0} r_\phi(s, 1/2) = r_\phi^*(1/2)$, and by construction $\Delta \geq 0$, so that its convex hull satisfies $\psi \geq 0$ and $\psi(0) = 0$. Thus

$$\begin{aligned} \psi(R(f) - R^*) &\leq \mathbb{E}[\psi(\mathbf{1}\{\text{sign}(f(X)) \neq \text{sign}(2\eta(X) - 1)\} |2\eta(X) - 1|)] \\ &= \mathbb{E}[\mathbf{1}\{\text{sign}(f(X)) \neq \text{sign}(2\eta(X) - 1)\} \psi(|2\eta(X) - 1|)] \end{aligned} \quad (16.1.6)$$

Now we use the special structure of the suboptimality function. By construction of ψ as the convex minorant of Δ , we have $\psi \leq \Delta$, and moreover, for any $s \in \mathbb{R}$

$$\begin{aligned} \mathbf{1}\{\text{sign}(s) \neq \text{sign}(2\eta - 1)\} \Delta(|2\eta - 1|) &= \mathbf{1}\{\text{sign}(s) \neq \text{sign}(2\eta - 1)\} \left[\inf_{s(2\eta-1) \leq 0} r_\phi(s, \eta) - r_\phi^*(\eta) \right] \\ &\leq r_\phi(s, \eta) - r_\phi^*(\eta), \end{aligned} \quad (16.1.7)$$

because $(1 + |2\eta - 1|)/2 = \max\{\eta, 1 - \eta\}$. Combining inequalities (16.1.6) and (16.1.7), we see that

$$\begin{aligned} \psi(R(f) - R^*) &\leq \mathbb{E}[\mathbf{1}\{\text{sign}(f(X)) \neq \text{sign}(2\eta(X) - 1)\} \Delta(|2\eta(X) - 1|)] \\ &\leq \mathbb{E}[r_\phi(f(X), \eta(X)) - r_\phi^*(\eta(X))] \\ &= R_\phi(f) - R_\phi^*, \end{aligned}$$

which is inequality (16.1.3).

The key to the claims (i)–(iii) in the second part of the theorem is the following lemma, which guarantees that so long as the sub-optimality function Δ_ϕ is positive, then so too is its biconjugate (16.1.2) ψ -transform.

Lemma 16.1.9. *The following hold:*

- (a) *The functions Δ and ψ are continuous on $[0, 1]$.*
- (b) *We have $\Delta \geq 0$ and $\Delta(0) = 0$.*
- (c) *If $\Delta(\delta) > 0$ for all $\delta > 0$, then $\psi(\delta) > 0$ for all $\delta > 0$.*

Deferring the proof of the lemma, We turn to the equivalence of items (i)–(iii) in the theorem. To see that classification calibration implies $\psi(\delta) \rightarrow 0$ if and only if $\delta \rightarrow 0$ (i.e., (i) implies (ii)), use Lemma 16.1.9: if ϕ is classification calibrated, then $\Delta(\delta) > 0$ for all $\delta > 0$ and by continuity, $\inf_{\delta \geq c} \psi(\delta) > 0$ for all $c > 0$, so that $\psi(\delta) \rightarrow 0$ if and only if $\delta \rightarrow 0$. Now, let (ii) hold: then if $R_\phi(f_n) \rightarrow R_\phi^*$, we necessarily have $\psi(R(f_n) - R^*) \rightarrow 0$ by inequality (16.1.3); this occurs if and only if $R(f_n) - R^* \rightarrow 0$ by assumption, that is, the implication (iii) holds. Finally, to see that if for any sequence of measurable functions f_n , the convergence $R_\phi(f_n) \rightarrow R_\phi^*$ implies $R(f_n) \rightarrow R^*$ (implication (iii)) the loss ϕ is necessarily classification calibrated (i). Assume for the sake of contradiction that (iii) holds but (i) fails, that is, ϕ is not classification calibrated. Then there must exist $\eta < 1/2$ and a sequence $s_n \geq 0$ (i.e. a sequence of predictions with incorrect sign) satisfying

$$r_\phi(s_n, \eta) \rightarrow r_\phi^*(\eta).$$

Construct the classification problem with a singleton $\mathcal{X} = \{x\}$, and set $\mathbb{P}(Y = 1) = \eta$. Then the sequence $f_n(x) = s_n$ satisfies $R_\phi(f_n) \rightarrow R_\phi^*$ but the true 0-1 risk $R(f_n) \not\rightarrow R^*$.

We return to the promised proof of Lemma 16.1.9.

Proof of Lemma 16.1.9 The function $r_\phi^*(\eta) = \inf_s \{\eta\phi(s) + (1 - \eta)\phi(-s)\}$ for $\eta \in [0, 1]$ and $r_\phi^*(\eta) = -\infty$ otherwise is a concave function, as it is the infimum of linear functions in η , and moreover, $-r_\phi^*$ is closed convex. Because ϕ is nonnegative, r_ϕ^* has domain $[0, 1]$. In one-dimension, closed convex functions are continuous on their domains (see Observation B.3.6 in Appendix B.3.2). Additionally, we also have

$$r_\phi^{\text{wrong}}(\eta) = \inf_{s \geq 0} \{\eta\phi(s) + (1 - \eta)\phi(-s)\}$$

for $0 \leq \eta < \frac{1}{2}$, and by symmetry $r_\phi^{\text{wrong}}(\eta) = r_\phi^{\text{wrong}}(-\eta)$, which is thus continuous on $[0, \frac{1}{2}) \cup (\frac{1}{2}, 1]$. A parallel argument to that for r_ϕ^* then shows that r_ϕ^{wrong} is left continuous and right continuous at $\frac{1}{2}$, and so

$$\Delta_\phi(\delta) = r_\phi^{\text{wrong}}\left(\frac{1 + \delta}{2}\right) - r_\phi^*\left(\frac{1 + \delta}{2}\right)$$

is continuous on $\delta \in [0, 1]$, where $\Delta_\phi(1 + \epsilon) = +\infty$. That $\psi = \Delta^{**}$ is continuous on $[0, 1]$ is then immediate by the continuity of closed convex functions on their domains.

The nonnegativity of Δ is immediate, and to see that $\Delta(0) = 0$, note that

$$r_\phi^{\text{wrong}}(1/2) := \inf_{s(1-1) \leq 0} r_\phi(s, 1/2) = \inf_s r_\phi(s, 1/2) = r_\phi^*(1/2),$$

so $\Delta(0) = r_\phi^*(1/2) - r_\phi^*(1/2) = 0$.

The final claim of the lemma that $\psi = \Delta^{**}$ is strictly positive whenever Δ is is more subtle. For this, we leverage the following technical lemma:

Lemma 16.1.10. *Let h be either (i) continuous on $[0, 1]$ or (ii) non-decreasing on $[0, 1]$, where $h(t) = +\infty$ for $t \notin [0, 1]$. If h satisfies $h(t) > 0$ for $t > 0$ and $h(0) = 0$, then $f(t) = h^{**}(t)$ satisfies $f(t) > 0$ for any $t > 0$.*

Proof We begin with the case (i). Define the function $h_{\text{low}}(t) := \inf_{s \geq t} h(s)$. Then because h is continuous, we know that over any compact set it attains its infimum, and thus (by assumption on h) $h_{\text{low}}(t) > 0$ for all $t > 0$. Moreover, h_{low} is non-decreasing. Define $f_{\text{low}}(t) = h_{\text{low}}^{**}(t)$ to be the biconjugate of h_{low} ; it is clear that $f \geq f_{\text{low}}$ as $h \geq h_{\text{low}}$. Thus we see that case (ii) implies case (i), so we turn to the more general result to see that $f_{\text{low}}(t) > 0$ for all $t > 0$.

For the result in case (ii), assume for the sake of contradiction there is some $z \in (0, 1)$ satisfying $h^{**}(z) = 0$. It is clear that $h^{**}(0) = 0$ and $h^{**} \geq 0$, so we must have $h^{**}(z/2) = 0$. Now, by assumption we have $h(z/2) = b > 0$, whence we have $h(1) \geq b > 0$. In particular, the piecewise linear function defined by

$$g(t) = \begin{cases} 0 & \text{if } t \leq z/2 \\ \frac{b}{1-z/2}(t - z/2) & \text{if } t > z/2 \end{cases}$$

is closed, convex, and satisfies $g \leq h$. But $g(z) > 0 = h^{**}(z)$, a contradiction to the fact that h^{**} is the largest (closed) convex function below h . \square

Lemma 16.1.10 immediately yields the final claim (c) of Lemma 16.1.9: we know that $\Delta(\delta) > 0$ and is continuous in δ . \square

JCD Comment: Put a figure

16.2 General surrogate risk consistency

The approach to proving consistency we develop for the binary classification cases extends, with some care, to essentially fully general scenarios. In this case, we consider a general supervised learning problem, where we assume that we have data in pairs $(X, Y) \in \mathcal{X} \times \mathcal{Y}$, where \mathcal{X} and \mathcal{Y} are general spaces. We assume we have a loss $\ell : \mathbb{R}^k \times \mathcal{Y} \rightarrow \mathbb{R}_+$ be a loss function we wish to minimize, so that the loss of a prediction function $f : \mathcal{X} \rightarrow \mathbb{R}^k$ for the pair (x, y) is $\ell(f(x), y)$. Let $\varphi : \mathbb{R}^k \times \mathcal{Y} \rightarrow \mathbb{R}$ be an arbitrary surrogate, where $\varphi(f(x), y)$ is the surrogate loss. Define the risk and φ -risk

$$R(f) := \mathbb{E}[\ell(f(X), Y)] \quad \text{and} \quad R_\varphi(f) := \mathbb{E}[\varphi(f(X), Y)].$$

Example 16.2.1 (Multiclass classification): Assume that $Y \in \{1, \dots, k\}$. For a score vector $s \in \mathbb{R}^k$, the zero-one error is

$$\ell_{0-1}(s, y) = \mathbf{1} \left\{ s_y \leq \max_{j \neq y} s_j \right\}.$$

Common surrogates include the multiclass logistic loss

$$\varphi(s, y) = \log \left(1 + \sum_{j \neq y} e^{s_j - s_y} \right)$$

and the pairwise comparison loss

$$\varphi(s, y) = \sum_{j \neq y} \phi(s_y - s_j),$$

where ϕ is a convex differentiable function with $\phi'(0) < 0$. These are both consistent losses (as we shall see). \diamond

Example 16.2.2 (Ranking problems): Suppose we wish to rank k items based on covariates x , so that $Y \in \mathfrak{S}_k$, the set of permutations of $[k]$, where we write $Y(i) \succ Y(j)$ to indicate that i is ranked ahead of j . A prediction assigning scores $s \in \mathbb{R}^k$ to each of the k items naturally induces a permutation via the sorting $s_{(1)} \geq s_{(2)} \geq \dots \geq s_{(k)}$. The Kendall tau distance counts the pairwise disagreements between s and Y via

$$\ell(s, y) = \binom{k}{2}^{-1} \sum_{i < j} \mathbf{1} \{s_i \geq s_j \text{ and } y(i) \prec y(j)\}.$$

In other scenarios, the data Y may consist of only partial feedback, such as a pairwise ranking, so that $Y = (i, j)$ indicates item i is preferred to j . Then for $y = (i, j)$, the pairwise disagreement becomes

$$\ell(s, y) = \mathbf{1} \{s_i \geq s_j\}.$$

These related losses—and others like them—admit essentially no efficiently minimizable surrogates. \diamond

Let $\mathcal{P}_{\mathcal{Y}}$ denote the space of all probability distributions on \mathcal{Y} , and define the conditional (pointwise) risks $r : \mathbb{R}^k \times \mathcal{P}_{\mathcal{Y}} \rightarrow \mathbb{R}$ and $r_{\varphi} : \mathbb{R}^k \times \mathcal{P}_{\mathcal{Y}} \rightarrow \mathbb{R}$ by

$$r(s, P) = \mathbb{E}_P[\ell(s, Y)] \quad \text{and} \quad r_{\varphi}(s, P) = \mathbb{E}_P[\varphi(s, Y)].$$

Following our development in the binary case, let $r^*(P) = \inf_s r(s, P)$ denote the minimal conditional risk, and similarly for $r_{\varphi}^*(P)$, when Y has distribution P . If $P_{Y|x}$ denotes the distribution of Y conditioned on $X = x$, then we may rewrite the risk functionals as

$$R(f) = \mathbb{E}[r(f(X), P_{Y|X})] \quad \text{and} \quad R_{\varphi}(f) = \mathbb{E}[r_{\varphi}(f(X), P_{Y|X})].$$

16.2.1 Uniform calibration

Continuing to follow the development in Section 16.1, for $\epsilon \geq 0$ define the *suboptimality gap function*

$$\Delta_{\varphi}(\epsilon, P) := \inf_{s \in \mathbb{R}^k} \{r_{\varphi}(s, P) - r_{\varphi}^*(P) \mid r(s, P) - r^*(P) \geq \epsilon\}, \quad (16.2.1)$$

which measures the gap between achievable (pointwise) risk and the best surrogate risk when we enforce that the true loss is not minimized to accuracy better than ϵ . We can then define the *uniform suboptimality function*

$$\Delta_{\varphi}(\epsilon) := \inf_{s \in \mathbb{R}^k, P \in \mathcal{P}_{\mathcal{Y}}} \{r_{\varphi}(s, P) - r_{\varphi}^*(P) \mid r(s, P) - r^*(P) \geq \epsilon\}. \quad (16.2.2)$$

(Exercise 16.6 shows how this relates to r_{ϕ} for binary classification.) Now let $\Delta_{\varphi}^{**}(\epsilon)$ be the biconjugate of Δ_{φ} , that is, Δ_{φ}^{**} is the largest convex function below Δ_{φ} . We can then make the following definition, which analogizes Definition 16.2.

Definition 16.3 (Uniform calibration). *The surrogate φ is uniformly calibrated for the loss ℓ if $\Delta_{\varphi}(\epsilon) > 0$ for all $\epsilon > 0$.*

We have the following proposition, which analogizes Theorem 16.1.5.

Proposition 16.2.3. *For any measurable $f : \mathcal{X} \rightarrow \mathbb{R}^k$,*

$$\Delta_{\varphi}^{**}(R(f) - R^*) \leq R_{\varphi}(f) - R_{\varphi}^*.$$

*Additionally, if φ is uniformly calibrated, then $\Delta_{\varphi}^{**}(\epsilon) > 0$ for all $\epsilon > 0$, and $R_{\varphi}(f_n) \rightarrow R_{\varphi}^*$ implies that $R(f_n) \rightarrow R^*$.*

Proof Let $\psi = \Delta_{\varphi}^{**}$ for shorthand. Then by Jensen's inequality,

$$\begin{aligned} \psi(R(f) - R^*) &= \psi(\mathbb{E}[r(f(X), P_{Y|X}) - r^*(P_{Y|X})]) \\ &\leq \mathbb{E}[\psi(r(f(X), P_{Y|X}) - r^*(P_{Y|X}))] \\ &\leq \mathbb{E}[\Delta_{\varphi}(r(f(X), P_{Y|X}) - r^*(P_{Y|X}))] \end{aligned}$$

where we use Jensen's inequality and that $\psi \leq \Delta_{\varphi}$. Now, note that by definition, for any $P \in \mathcal{P}_{\mathcal{Y}}$ and $f(x) \in \mathbb{R}^k$, we have

$$\begin{aligned} \Delta_{\varphi}(r(f(x), P) - r^*(P)) &= \inf_{s \in \mathbb{R}^k} \{r_{\varphi}(s, P) - r_{\varphi}^*(P) \mid r(s, P) - r^*(P) \geq r(f(x), P) - r^*(P)\} \\ &\leq r(f(x), P) - r^*(P). \end{aligned}$$

Substituting into the preceding display, we obtain

$$\psi(R(f) - R^*) \leq \mathbb{E} [r_\varphi(f(X), P_{Y|X}) - r_\varphi^*(P_{Y|X})] = R_\varphi(f) - R_\varphi^*,$$

giving the first statement of the proposition. For the second statement, note that $\epsilon \mapsto \Delta_\varphi(\epsilon)$ is non-decreasing by construction, and uniform calibration implies it is strictly positive for $\epsilon > 0$. Lemma 16.1.10 implies its biconjugate is positive as well. \square

16.2.2 Pointwise calibration

Unfortunately, it can be challenging to explicitly verify uniform calibration, as it involves several nested infima over various potentially non-convex sets. Thus, it is frequently convenient to relax to loss calibration:

Definition 16.4. *The surrogate φ is calibrated for the loss ℓ if for all $P \in \mathcal{P}_Y$ and $\epsilon > 0$,*

$$\Delta_\varphi(\epsilon, P) > 0.$$

Definition 16.4 relaxes the uniform calibration condition in Definition 16.3 to apply pointwise, that is, for each distribution P on \mathcal{Y} . While this precludes the cleanest guarantee that $\psi(R_\varphi(f) - R_\varphi^*) \leq R(f) - R^*$, under a minor condition to address integrability issues (see also Exercise 16.7), it still is sufficient to guarantee surrogate risk consistency.

Theorem 16.2.4. *Let ℓ be bounded. Then φ is calibrated for the loss ℓ if and only if it is surrogate risk consistent for ℓ , that is, $R_\varphi(f_n) \rightarrow R_\varphi^*$ implies that $R(f_n) \rightarrow R^*$.*

Proof We prove the implication that φ is calibrated implies that it is surrogate risk consistent; the converse is an essentially immediate exercise by considering distributions supported on a single point x . Let $B < \infty$ be the bound on ℓ , so we may assume that $r(s, P) - r^*(P) \leq B$. Let f_n be a sequence of functions satisfying $R_\varphi(f_n) \rightarrow R_\varphi^*$ and let $\epsilon > 0$ be arbitrary. Let $P_x = P(Y \in \cdot \mid X = x)$ be the conditional distribution of Y given $X = x$, and for shorthand, define

$$\delta_n(x) := r(f_n(x), P_x) - r^*(P_x) \quad \text{and} \quad A_{n,\epsilon} := \{x \in \mathcal{X} \mid \delta_n(x) \geq \epsilon\}.$$

so that $\delta_n(x) \in [0, B]$. Then for the risk gap $R(f_n) - R^*$, we see that

$$R(f_n) - R^* = \mathbb{E}[\delta_n(X)] = \mathbb{E}[\mathbf{1}\{X \in A_{n,\epsilon}\} \delta_n(X)] + \underbrace{\mathbb{E}[\mathbf{1}\{X \notin A_{n,\epsilon}\} \delta_n(X)]}_{\leq \epsilon},$$

where we have used that $\delta_n(x) \leq \epsilon$ on $A_{n,\epsilon}^c$ by definition. Now, consider the first expectation. Let $\gamma > 0$ be otherwise arbitrary, and note that

$$\begin{aligned} \mathbf{1}\{x \in A_{n,\epsilon}\} &= \mathbf{1}\{\delta_n(x) \geq \epsilon\} \leq \frac{r_\varphi(f_n(x), P_x) - r_\varphi^*(P_x)}{\Delta_\varphi(\epsilon, P_x)} \mathbf{1}\{\Delta_\varphi(\epsilon, P_x) > \gamma\} + \mathbf{1}\{\Delta_\varphi(\epsilon, P_x) \leq \gamma\} \\ &\leq \frac{r_\varphi(f_n(x), P_x) - r_\varphi^*(P_x)}{\gamma} + \mathbf{1}\{\Delta_\varphi(\epsilon, P_x) \leq \gamma\}. \end{aligned}$$

In particular, as $\delta_n(x) \in [0, B]$, we have

$$\begin{aligned} \mathbb{E}[\delta_n(X)\mathbf{1}\{X \in A_{n,\epsilon}\}] &\leq B\mathbb{E}[\mathbf{1}\{X \in A_{n,\epsilon}\}] \\ &\leq B \left[\frac{\mathbb{E}[r_\varphi(f_n(X), P_X) - r_\varphi^*(P_X)]}{\gamma} + \mathbb{P}(\Delta_\varphi(\epsilon, P_X) \leq \gamma) \right] \\ &= B \left[\frac{R_\varphi(f_n) - R_\varphi^*}{\gamma} + \mathbb{P}(\Delta_\varphi(\epsilon, P_X) \leq \gamma) \right] \end{aligned}$$

for all $\gamma \geq 0$. But of course, by assumption the first term satisfies $R_\varphi(f_n) - R_\varphi^* \rightarrow 0$, while the continuity of probability measures guarantees that $\lim_{\gamma \downarrow 0} \mathbb{P}(\Delta_\varphi(\epsilon, P_X) \leq \gamma) = \mathbb{P}(\Delta_\varphi(\epsilon, P_X) = 0) = 0$. In particular, $\mathbb{E}[\mathbf{1}\{X \in A_{n,\epsilon}\} \delta_n(X)] \rightarrow 0$, and so $R(f_n) - R^* \leq \epsilon + o(1)$ as $n \rightarrow \infty$. As $\epsilon > 0$ was arbitrary, this completes the proof. \square

16.2.3 Examples: multiclass surrogate risk consistency

Happily, it is much easier to provide examples of the pointwise calibration condition in Definition 16.4 than the uniform calibration guarantee. Here, we focus on multiclass calibration guarantees, where $Y \in \{1, \dots, k\}$ represents one of k classes, and we consider predictors $f : \mathcal{X} \rightarrow \mathbb{R}^k$ and consistency for the zero-one loss. While no formulation of surrogate risk consistency as clean as that Theorem 16.1.5 provides for the margin-based losses for binary classification exists, we can still give clean sufficient conditions for surrogate risk consistency. Here, we identify distributions on Y with vectors $p \in \Delta_k = \{p \in \mathbb{R}_+^k \mid \langle \mathbf{1}, p \rangle = 1\}$, so that $r(s, p) = \sum_{y=1}^k p_y \mathbf{1}\{s_y \leq \max_{j \neq y} s_j\}$.

Let us revisit Example 16.2.1 to make this more concrete. We first consider pairwise margin-based losses of the form

$$\varphi(s, y) = \sum_{j=1}^k \phi(s_y - s_j), \quad (16.2.3)$$

where $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is a non-increasing function. The intuition is similar to the classification calibration case: to minimize this, we expect $s_y \gg s_j$ for $j \neq y$, as ϕ is non-increasing. To make this more precise, we have the following lemma.

Lemma 16.2.5. *Let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ be non-increasing and satisfy $\phi(t) < \phi(-t)$ for all $t > 0$ and the surrogate φ take the pairwise form (16.2.3). Let $p \in \Delta_k$ and $s \in \mathbb{R}^k$ satisfy $r_\varphi(s, p) = r_\varphi^*(p)$. If $p_i > p_j$, then $s_i \geq s_j$, and if ϕ is differentiable with $\phi'(0) < 0$, then $s_i < s_j$.*

Proof Without loss of generality, we take $i = 1$ and $j = 2$, so that $p_1 > p_2$. Let $s \in \mathbb{R}^k$ and $s' = (s_2, s_1, s_3^k)$ be s with its first two entries flipped. Then

$$\begin{aligned} r_\varphi(s, p) - r_\varphi(s', p) &= p_1 \phi(s_1 - s_2) + p_1 \sum_{j>3} \phi(s_1 - s_j) + p_2 \phi(s_2 - s_1) + p_2 \sum_{j>3} \phi(s_2 - s_j) \\ &\quad - \left(p_2 \phi(s_1 - s_2) + p_2 \sum_{j>3} \phi(s_1 - s_j) + p_1 \phi(s_2 - s_1) + p_1 \sum_{j>3} \phi(s_2 - s_j) \right) \\ &= (p_1 - p_2) \left(\phi(s_1 - s_2) - \phi(s_2 - s_1) + \sum_{j>3} (\phi(s_1 - s_j) - \phi(s_2 - s_j)) \right). \end{aligned}$$

If $s_1 < s_2$, then $\phi(s_1 - s_j) \geq \phi(s_2 - s_j)$ and $\phi(s_1 - s_2) > \phi(s_2 - s_1)$, by assumption, so that if $p_1 > p_2$, we must have $r_\varphi(s, p) - r_\varphi(s', p) > (p_1 - p_2) \cdot (0 + 0) = 0$, which would contradict that $r_\varphi(s, p) = r_\varphi^*(p)$.

When ϕ is differentiable, if s minimizes $r_\varphi(s, p)$, then we have $\nabla r_\varphi(s, p) = 0$, and in particular,

$$0 = p_1 \sum_{j=1}^k \phi'(s_1 - s_j) - \sum_{j=1}^k p_j \phi'(s_j - s_1).$$

Assume for the sake of contradiction that $s_1 = s_2 = t$ for some $t \in \mathbb{R}$. Then the preceding equality implies that

$$p_1 \sum_{j=1}^k \phi'(t - s_j) = p_2 \sum_{j=1}^k \phi'(t - s_j), \quad \text{i.e. } p_1 \phi'(0) = p_2 \phi'(0).$$

But as $p_1 > p_2$ and $\phi'(0) < 0$, this necessarily fails. \square

This lemma implies the following proposition, which gives sufficient conditions for a convex margin-based loss to be surrogate risk consistent.

Proposition 16.2.6. *Let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ be non-increasing, convex, and differentiable with $\phi'(0) < 0$, and let the surrogate φ take the pairwise form (16.2.3). Then φ is surrogate risk consistent.*

The differentiability assumption on ϕ is important here; without it, Proposition 16.2.6 may fail (see Exercise 16.8).

Proof We use Theorem 16.2.4. Fix $p \in \Delta_k$ and without loss of generality assume that $p_1 > p_2$; by Theorem 16.2.4, it is then enough to show that

$$\inf_s \{r_\varphi(s, p) - r_\varphi^*(p) \mid s_1 \leq s_2\} > 0.$$

Let $s^{(n)} \in \mathbb{R}^k$ be any sequence with $s_1^{(n)} \leq s_2^{(n)}$ for all n and assume for the sake of contradiction that $r_\varphi(s^{(n)}, p) \rightarrow r_\varphi^*(p)$. Then (working with the compactification $[-\infty, \infty]$ of \mathbb{R} as necessary) there is a convergent subsequence, which w.l.o.g. we may take as the entire sequence, with $s^{(n)} \rightarrow s \in \{\mathbb{R} \cup \{\pm\infty\}\}^k$, where $s_1 \leq s_2$. But then (with trivial modification to address the infinite case) this contradicts Lemma 16.2.5. \square

One can demonstrate other variants of Proposition 16.2.6 using similar loss functions. For example, we have the following, whose proof Exercise 16.9 outlines.

Proposition 16.2.7. *Let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ be non-increasing, convex, and differentiable with $\phi'(0) < 0$. Let the surrogate*

$$\varphi(s, y) := \phi(s_y) + \sum_{j \neq y} \phi(-s_j).$$

Then φ is surrogate risk consistent.

16.3 Generalized entropies and surrogate risk consistency

In Section 14.3.1 of Chapter 14, we develop a theory of exact equivalence between a proper loss ℓ and a convex surrogate φ , though with a slightly different perspective than what we take here, as the goal is to prediction values $\mu \in \mathcal{M} = \{\mathbb{E}_P[Y] \mid P \in \mathcal{P}_Y\}$. We can provide a parallel development here, showing that for general losses on vectors $s \in \mathbb{R}^k$, we have a natural generalized entropy notion, and from any generalized entropy, we can construct a surrogate loss that (often) is guaranteed to be surrogate risk consistent. Thus, entropies once again provide an explicit link between loss minimization and consistency, though we shall weaken the requirement that the two entropies a loss and its surrogate generate be equal; equality will allow us to develop more nuanced notions of consistency that we explore in Section 16.5 to come.

We focus first on the case of multiclass classification. Recall that a (generalized) entropy $H : \Delta_k \rightarrow \mathbb{R}$ is any concave function H with $H(p) > -\infty$ except (potentially) when $p_j = 0$ for some j . We follow the notation of Chapter 14 and let

$$\Omega(p) = -H(p)$$

be the negative entropy, and require that H be closed, meaning that Ω is a closed convex function. We immediately see that

$$\varphi(s, y) := -s_y + \Omega^*(s), \quad (16.3.1)$$

where

$$\Omega^*(s) = \sup_{p \in \Delta_k} \{\langle p, s \rangle + H(p)\}$$

is the convex conjugate of $\Omega = -H$, satisfies

$$\inf_s \mathbb{E}_p[\varphi(s, Y)] = \inf_s \{\Omega^*(s) - \langle p, s \rangle\} = -\Omega^{**}(p) = -\Omega(p) = H(p).$$

We can record this as a proposition, where we recall the entropy $H_\ell(p) := \inf_s \mathbb{E}_p[\ell(s, Y)]$ associated to the loss ℓ from Chapter 14.

Proposition 16.3.1. *Let $H : \Delta_k \rightarrow \mathbb{R}$ be closed concave and $\Omega = -H$. Then the losses $\varphi(s, y) = -s_y + \Omega^*(s)$ are closed, convex, and satisfy*

$$H_\varphi(p) := \inf_{s \in \mathbb{R}^k} \mathbb{E}_p[\varphi(s, Y)] = H(p).$$

So any (generalized) entropy on the simplex *generates* a convex loss with the same entropy.

JCD Comment: Perhaps reference earlier material a little bit more carefully here.

The construction (16.3.1) is a privileged construction, as from it we can derive desirable properties of the convex loss φ itself, and in particular, the loss φ is frequently surrogate risk consistent (Definition 16.1) for the zero-one error. Specializing Definition 16.4 to the classification case, we consider the following definition.

Definition 16.5. *The loss $\varphi : \mathbb{R}^k \times [k] \rightarrow \mathbb{R}$ is classification calibrated for the zero-one loss if for any $p \in \Delta_k$ and y such that $p_y < \max_j p_j$,*

$$\inf_s \left\{ \sum_{j=1}^k p_j \varphi(s, j) \right\} < \inf_s \left\{ \sum_{j=1}^k p_j \varphi(s, j) \mid s_y \geq \max_j s_j \right\}.$$

Theorem 16.2.4 shows that Definition 16.5 is equivalent to surrogate risk consistency.

The generalized entropy function H provides a convenient route to determining the classification calibration of the surrogate φ from the construction (16.3.1). To provide sufficient conditions, we recall the collection of uniformly convex functions (see Appendix C.1.2, equation (C.1.5)):

Definition 16.6. A convex function $f : \mathbb{R}^k \rightarrow \overline{\mathbb{R}}$ is $(\lambda, \kappa, \|\cdot\|)$ -uniformly convex over $C \subset \mathbb{R}^k$ if it is closed and for all $t \in [0, 1]$ and $x_1, x_2 \in C$,

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2) - \frac{\lambda}{2}t(1-t)\|x_1 - x_2\|^\kappa [(1-t)^{\kappa-1} + t^{\kappa-1}].$$

Uniform convexity is a weaker version of strong convexity (which corresponds to $\kappa = 2$), and guarantees that the function exhibits superlinear growth. The condition is sufficient to guarantee the classification calibration of the associated loss φ :

Theorem 16.3.2. Let H be closed concave, symmetric, and have domain $\text{dom } H = \Delta_k$. Let φ have the definition (16.3.1). Assume that (a) H is strictly concave and $\inf_s \sum_{j=1}^k p_j \ell(s, j)$ is attained for each $p \in \Delta_k$, or (b) H is uniformly concave. Then φ is classification calibrated.

As should not be surprising, this result relies strongly on the stability properties of the minimizers of the loss φ as well as the differentiability properties of the conjugate Ω^* of the negative entropy $\Omega = -H$. We therefore postpone the proof to Section 16.3.1, instead using it to give a few more concrete examples of consistent losses.

Example 16.3.3 (Multiclass logistic regression): The multiclass logistic loss

$$\varphi(s, y) = \log \left(\sum_{j=1}^k \exp(s_j - s_y) \right) = -s_y + \log \left(\sum_{j=1}^k e^{s_j} \right)$$

has associated entropy

$$H(p) = \inf_s \mathbb{E}_p[\varphi(s, Y)] = - \sum_{j=1}^k p_j \log p_j,$$

the classical entropy of Y . (Recall Examples 14.3.4 and 14.3.6.) This is strongly convex over Δ_k by Pinsker's inequality, and the negative entropy $\Omega = -H$ has the dual

$$\Omega^*(s) = \log \left(\sum_{j=1}^k e^{s_j} \right),$$

showing that the multiclass logistic loss is consistent. \diamond

Example 16.3.4 (Norm-based losses): For any $1 < \kappa < \infty$, take any loss of the form

$$\varphi(s, y) = -s_y + \frac{1}{\kappa} \|s\|_\kappa^\kappa,$$

which for the conjugate $\kappa^* = \frac{\kappa}{\kappa-1} < \infty$ satisfies

$$H(p) = \inf_s \mathbb{E}_p[\varphi(s, Y)] = -\frac{1}{\kappa^*} \|p\|_{\kappa^*}^{\kappa^*}.$$

This entropy is strictly concave, and the infimum is always attained, so Theorem 16.3.2 gives that φ is surrogate-risk consistent. \diamond

16.3.1 Proof of Theorem 16.3.2

We will freely use the conjugate duality relationships we provide in Appendices C.1.2 and C.2.2. Let $\Omega = -H$ be the (convex) negative entropy, noting that $\Omega^*(s) = \sup_{p \in \Delta_k} \{\langle p, s \rangle + H(s)\} < \infty$, as suprema of closed concave functions over compact sets are attained, and thus $\text{dom } \Omega^* = \mathbb{R}^k$. Note also that Ω^* is continuously differentiable under the conditions of the theorem, as Ω is strictly convex, because uniform convexity is stronger than strict convexity. (See Propositions C.2.7 and C.2.8.)

We begin with an intermediate lemma.

Lemma 16.3.5. *Under the conditions of Theorem 16.3.2, Ω^* is continuously differentiable on \mathbb{R}^k , satisfies $\Delta_k \subset \nabla \Omega^*(\mathbb{R}^k)$, and if $s_i \geq s_j$, then $p = \nabla \Omega^*(s)$ satisfies $p_i \geq p_j$.*

Proof The first claims of the lemma we have already demonstrated; all that remains is to show that if $s_i \geq s_j$ then $p_i \geq p_j$. We know that $p = \nabla \Omega^*(s) = \text{argmax}_{p \in \Delta_k} \{\langle p, s \rangle - \Omega(p)\}$, and let us assume for the sake of contradiction that $p_i < p_j$. Then letting p' be the vector p with entries i and j swapped, we have

$$\Omega(p') = \Omega(p) \quad \text{but} \quad \Omega\left(\frac{1}{2}(p' + p)\right) < \frac{1}{2}\Omega(p') + \frac{1}{2}\Omega(p) = \Omega(p') = \Omega(p).$$

We also have

$$\langle s, p \rangle - \langle s, p' \rangle = (s_i - s_j)(p_i - p_j) \leq 0,$$

that is, $\langle s, p \rangle \leq \langle s, p' \rangle$. But then

$$\frac{1}{2}\langle s, p + p' \rangle - \Omega\left(\frac{1}{2}p + \frac{1}{2}p'\right) \geq \langle s, p \rangle - \Omega\left(\frac{1}{2}p + \frac{1}{2}p'\right) > \langle s, p \rangle - \Omega(p),$$

a contradiction to the assumed optimality that $p = \nabla \Omega^*(s)$. \square

From this lemma, we obtain the following immediate observation:

Observation 16.3.6. *Under the conditions of Theorem 16.3.2, if $s^* \in \mathbb{R}^k$ minimizes $\mathbb{E}_p[\varphi(s, Y)]$ then $p_i > p_j$ implies that $s_i^* > s_j^*$.*

Proof We know that s^* minimizes $\mathbb{E}_p[\varphi(s, Y)]$ if and only if $p = \nabla \Omega^*(s^*)$ by the Fenchel-Young inequality (14.1.4). Take the contrapositive of Lemma 16.3.5. \square

We also require the following technical lemma, whose proof we defer temporarily:

Lemma 16.3.7. *Let f have Hölder-continuous gradient and assume $f^* = \inf_u f(u) > -\infty$. Then if $f(u_n) \rightarrow f^*$, we have $\nabla f(u_n) \rightarrow 0$.*

Now we may prove the theorem. Under assumption (a) of the theorem, that H is strictly concave and that the infimum $\inf_s \mathbb{E}_p[\varphi(s, Y)]$ is always attained, the preceding observation immediately gives the result: if $p_i > p_j$, then we must have

$$\inf_s \{\mathbb{E}_p[\varphi(s, Y)] \mid s_i \leq s_j\} > \inf_s \mathbb{E}_p[\varphi(s, Y)],$$

which is classification calibration (Definition 16.5). Under assumption (b), that H is uniformly convex, recall that $H(p) = \inf_s \{\sum_{j=1}^k p_j \varphi(s, j)\} > -\infty$, and let $p_i > p_j$. Let $s^{(m)}$ be any sequence

such that $\sum_{j=1}^k p_j \varphi(s^{(m)}, j) \rightarrow H(p)$. Because H is uniformly concave (i.e., Ω is uniformly convex), the dual gradient $\nabla \Omega^*$ is Hölder continuous with $\text{dom } \nabla \Omega^* = \mathbb{R}^k$ (Proposition C.2.7). Then because $\sum_{j=1}^k p_j \varphi(s, j) = -\langle p, s \rangle + \Omega^*(s)$, we obtain

$$\lim_{m \rightarrow \infty} \nabla \Omega^*(s^{(m)}) = p$$

by Lemma 16.3.7. If $s_i^{(m)} \geq s_j^{(m)}$, then Lemma 16.3.5 implies that $p^{(m)} = \nabla \Omega^*(s^{(m)})$ satisfies $p_i^{(m)} \geq p_j^{(m)}$, and so if $p_i > p_j$, it must be the case that eventually $s_i^{(m)} > s_j^{(m)}$. Moreover, continuity of $\nabla \Omega^*$ shows that

$$\liminf_m |s_i^{(m)} - s_j^{(m)}| = 0 \text{ implies } \liminf_m |p_i^{(m)} - p_j^{(m)}| = 0,$$

and so if $p_i > p_j$, then we necessarily have $\liminf_m (s_i^{(m)} - s_j^{(m)}) > 0$. In particular, we have demonstrated that

$$\inf_s \left\{ \sum_{j=1}^k p_j \varphi(s, j) \mid s_i \leq s_j \right\} > H(p)$$

whenever $p_i > p_j$, which is classification calibration (Definition 16.5).

We finally return to the promised proof of Lemma 16.3.7.

Proof of Lemma 16.3.7 Let $\beta > 0$ be the Hölder constant for the ℓ_2 -norm, so that for $\|\cdot\| = \|\cdot\|_2$ there exists $C < \infty$ such that $\|\nabla f(u) - \nabla f(v)\| \leq C \|u - v\|^\beta$ for all u, v . Then using that $h(t) = f(u + t(v - u))$ satisfies $h'(t) = \langle \nabla f(u + t(v - u)), v - u \rangle$ and $f(v) = h(1) = h(0) + \int_0^1 h'(t) dt$, we have

$$\begin{aligned} f(v) &= f(u) + \langle \nabla f(u), v - u \rangle + \int_0^1 \langle \nabla f(u + t(v - u)) - \nabla f(u), v - u \rangle dt \\ &\leq f(u) + \langle \nabla f(u), v - u \rangle + \int_0^1 C t^\beta \|u - v\|^{\beta+1} dt = f(u) + \langle \nabla f(u), v - u \rangle + \frac{C}{\beta+1} \|u - v\|^{\beta+1}. \end{aligned}$$

Now, for the sake of contradiction, let u_n satisfy $f(u_n) \rightarrow \inf_u f(u)$ and $g_n = \nabla f(u_n) \not\rightarrow 0$. Then there exists some $c > 0$ such that $\|g_n\| \geq c$ infinitely often; without loss of generality, assume this is the full sequence. Then fixing $\delta > 0$ to be chosen, let $v_n = u_n - \delta g_n / \|g_n\|$, so that

$$f(v_n) \leq f(u_n) - \delta \|g_n\| + \frac{C}{\beta+1} \delta^{\beta+1} \leq f(u_n) + \delta \left(\frac{C}{\beta+1} \delta^\beta - c \right).$$

Take any δ small enough that the last quantity is strictly negative to obtain that $\limsup_n f(v_n) < \liminf_n f(u_n) = f^*$, a contradiction. \square

16.4 Structured prediction and generalized entropies

Structured prediction problems involve predicting objects more complicated than standard classification problems, and their development allows us to apply statistical machine learning techniques to more sophisticated prediction problems. In Example 16.2.2, for instance, we wished to predict rankings of k items, and more broadly, we frequently wish to predict some type of combinatorial object $y \in \mathcal{Y}$, such as the matching in a graph, the alignment of images, a tree indicating some dependence structure, among other possibilities.

Example 16.4.1 (Matching): In protein structure prediction, we are given a sequence of amino acids and wish to determine the 3-dimensional structure of the resulting molecule. Often, cysteine (one of the amino acids) molecules bind to one another in pairs, offering clues to the resulting structure. Predicting these bindings is then equivalent to predicting a matching in a graph: given k cysteine molecules, we determine which $k/2$ pairs are matched. \diamond

Example 16.4.2 (Image registration): The *image registration* problem arises when one is given a sequence of images of a particular environment and wishes to identify objects within sequential images. We can represent this as follows: consider two $d \times d$ images, collections of pixels $x \in [0, 1]^{d \times d}$ and $x' \in [0, 1]^{d \times d}$ (for simplicity, we assume the images only have a pixel intensity and so are grayscale, though multi-channel images add little additional complexity). The registration problem is to match pixels x_{ij} in the first image to pixels $x'_{i'j'}$ in the second. Assuming it is possible for pixels to be unmatched, a registration consists of a mapping $\pi : [d] \times [d] \rightarrow ([d] \times [d] \cup \{\emptyset\})$, where $\pi(i, j) = (i', j')$ indicates the pixel at location (i, j) in image x matches to (i', j') in image x' (and $\pi(i, j) = \emptyset$ that the pixel is unmatched). \diamond

Abstracting away the details in the two examples, the typical setting we consider is that $y \in \mathcal{Y}$ consists of “parts,” where each part contributes to the whole prediction. In Example 16.4.1, parts are pairs of nodes, each with a particular binding affinity; in Example 16.4.2, we may represent parts as pairs of nodes in a bipartite graph. So in the former, given a graph on k nodes, a prediction is a subset $y \subset [k] \times [k]$ of edges, where for perfect matchings, each node i matches with exactly one other node j (and not itself), and for symmetry we can represent $(i, j) \in y$ if and only if $(j, i) \in y$. In the latter, we are given two sets of k items, and wish to find a matching between the two; equivalently, a permutation of $[k]$, which we can therefore represent as $y \in \mathfrak{S}_k$ the set of permutations of $[k]$, where $y(i)$ represents the index to which i matches. Figure 16.1 represents this graphically. In each of these cases, it is natural to represent the loss of a predicted matching y' for a correct matching y by counting the number of edges incorrectly labeled.

In structured prediction settings, we typically assume the loss of interest decomposes across the “parts” of y , so that each part of y contributes to the loss as a whole. We will consequently consider a broad family of structured losses that we can write as follows. We have a statistic $\tau : \mathcal{Y} \rightarrow \mathbb{R}^m$ representing the parts of y , and the loss of a prediction $y' \in \mathcal{Y}$ for a correct value $y \in \mathcal{Y}$ is

$$\ell(y', y) = \tau(y')^\top A \tau(y) + c \quad (16.4.1)$$

where $A \in \mathbb{R}^{m \times m}$ is a matrix and c is a scalar, and we assume

$$y = \operatorname{argmin}_{y' \in \mathcal{Y}} \ell(y', y) \quad (16.4.2)$$

To alleviate the abstractness of the formulation (16.4.1), we rewrite the problems we have already considered in its form.

Example 16.4.3 (Cost-weighted multiclass classification): Let the label $y \in \{1, \dots, k\}$, and let $w_{ij} \geq 0$ be the loss for predicting class i when the true class is j (and $w_{ii} = 0$). Then we take $\tau(y) = e_y$, the y th standard basis vector, and $A \in \mathbb{R}^{k \times k}$ with entries $A_{ij} = w_{ij}$, which gives $\ell(y', y) = \tau(y')^\top A \tau(y) = w_{y', y}$ in the formulation (16.4.1). \diamond

Example 16.4.4 (Ranking and bipartite matchings): For bipartite matching, which is equivalent to predicting a permutation, we let $y, y' \in \mathfrak{S}_k$ be permutations. Consider the Hamming

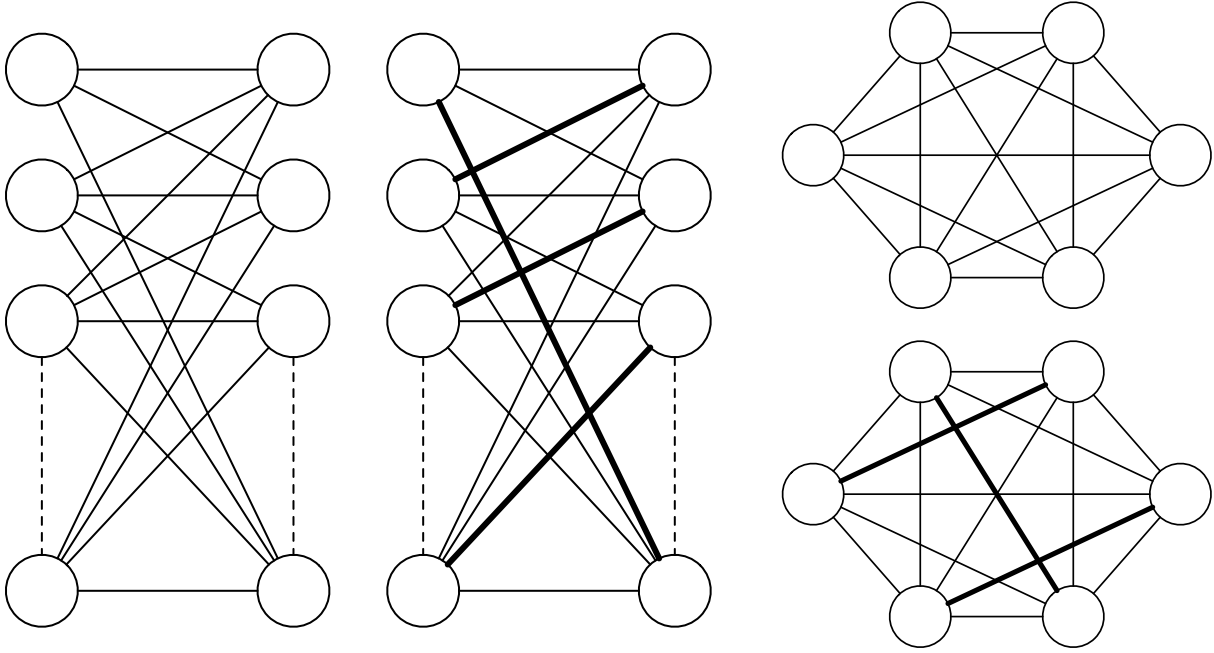


Figure 16.1. Left two figures: a bipartite graph, with edges across, and a bipartite graph with matching indicated. Right two figures: a fully connected graph and fully connected graph exhibiting a perfect matching, with each node matched to exactly one other node.

distance between permutations $\ell(y', y) = \sum_{i=1}^k \mathbf{1}\{y'(i) \neq y(i)\}$. Then for $\langle \cdot, \cdot \rangle$ the usual inner product on matrices, we let $\tau(y) \in \{0, 1\}^{k \times k}$ be the matrix representing the permutation y , that is, $M = \tau(y)$ satisfies $M_{i, y(i)} = 1$ for each i and zeros elsewhere. Then

$$\ell(y', y) = k - \langle \tau(y'), \tau(y) \rangle = k - \sum_{i=1}^k \mathbf{1}\{y'(i) = y(i)\}$$

has the representation (16.4.1). Alternatively, Kendall's *tau distance* counts the pairwise disagreements between permutations $y, y' : [k] \rightarrow [k]$ via

$$\ell(y', y) = \sum_{1 \leq i, j \leq k} \mathbf{1}\{y'(i) < y'(j) \text{ and } y(i) > y(j)\} + \mathbf{1}\{y'(i) > y'(j) \text{ and } y(i) < y(j)\}.$$

Recognizing this as the sum of discordant pairs and letting $\text{sign}(0) = 0$, we rewrite

$$\ell(y', y) = \sum_{1 \leq i, j \leq k} (1 - \text{sign}(y'(i) - y'(j)) \text{sign}(y(i) - y(j))),$$

so we see that if we define the matrix $\tau(y) \in \{-1, 0, 1\}^{k \times k}$ by $\tau(y)_{ij} = \text{sign}(y(i) - y(j))$, then $\ell(y', y) = k^2 - \langle \tau(y'), \tau(y) \rangle$, which again is of the form (16.4.1). \diamond

Example 16.4.5 (Perfect matchings): For predicting a (perfect) matching in a graph of k nodes, we identify a matching y as a collection of edges $(i, j) \in [k] \times [k]$ (where we take the convention that $(i, j) \in y$ if and only if $(j, i) \in y$ to address symmetry). Then the loss

$\ell(y', y) = \sum_{(i,j) \in y'} \mathbf{1}\{(i,j) \notin y\}$ counts the number of edges present in y' but not in y (which, by symmetry, is the number of edges in y not in y'). Rewriting,

$$\ell(y', y) = \frac{k}{2} - \text{card}(y \cap y') = \frac{k}{2} - \frac{1}{2} \sum_{i,j=1}^k \mathbf{1}\{(i,j) \in y'\} \mathbf{1}\{(i,j) \in y\}.$$

Letting $\langle \cdot, \cdot \rangle$ be the usual inner product on matrices and $\tau(y) \in \mathbb{R}^{k \times k}$ have entries $\tau(y)_{ij} = \mathbf{1}\{(i,j) \in y\}$, we have $\ell(y', y) = -\frac{1}{2} \langle \tau(y'), \tau(y) \rangle + \frac{k}{2}$, agreeing with the formulation (16.4.1). \diamond

Given a loss with the representation (16.4.1), the the prediction function $\hat{y}(s)$ from a vector $s \in \mathbb{R}^m$ of scores we use is

$$\hat{y}(s) \equiv \text{pred}(s) := \operatorname{argmax}_{y \in \mathcal{Y}} \left\{ \tau(y)^\top s \right\},$$

which decodes a score vector s into a “most likely” \hat{y} ; we choose the element arbitrarily if the maximizer is non-unique. It will turn out that the generalized entropy associated to the loss (16.4.1) and an analogue of the surrogate (16.3.1) will be consistent for predictions here, showing that the construction (16.3.1) is indeed a privileged one, providing a natural approach (with guaranteed surrogate-risk consistency properties) for constructing convex loss functions.

16.4.1 The failure of naive margin- and hinge-type losses

Before developing the entropy-based surrogate losses for general structured prediction problems, we take a brief historical detour. Because the problem of learning a structured predictor is to finding scores s so that the prediction $\hat{y}(s) = \operatorname{argmax}_y \{\tau(y)^\top s\}$ is typically correct, and we wish to have a convex loss function, a natural first attempt at losses is to generalize the margin-based losses from Section 16.1. For this, we require a notion of the “margin” of a prediction. Using the prediction scheme $\hat{y}(s)$, we would like to find a vector $s \in \mathbb{R}^k$ so that for a true label y , we assign values $\tau(y')^\top s \ll \tau(y)^\top s$ for y' far from y , and to reflect that we have a loss ℓ , we will seek $\tau(y')^\top s \leq \tau(y)^\top s - \ell(y, y')$. Then the most natural notion of a margin is

$$\min_{\hat{y} \in \mathcal{Y}} \left\{ s^\top (\tau(y) - \tau(\hat{y})) - \ell(\hat{y}, y) \right\},$$

which is positive so long as $s^\top \tau(y) \geq s^\top \tau(y') + \ell(y', y)$ for all $y' \in \mathcal{Y}$.

In analogy with the hinge loss (recall Examples 16.1.1 and 16.1.3), then the standard choice in structured prediction problems is to choose losses to maximize the margin, leading to surrogates of the form

$$\phi(s, y) := \max_{\hat{y} \in \mathcal{Y}} \left\{ \ell(\hat{y}, y) + s^\top (\tau(\hat{y}) - \tau(y)) \right\}, \quad (16.4.3)$$

the *structured prediction hinge loss*. (See, for example, [58, 122, 178, 177].) Then so long as $\ell(y', y) > 0$ for $y' \neq y$, it is immediate that $\phi(s, y) = 0$ implies $\hat{y}(s) = y$ uniquely; moreover, the loss (16.4.3) is convex as it is the maximum of linear functions of s .

While we will not dwell on this, the formulation (16.4.3) also has the major advantage that for many structured sets \mathcal{Y} —even exponentially large sets—computing the maximum is computationally efficient. For example, in Examples 16.4.4 and 16.4.5, computing the objective (16.4.3) is equivalent to solving a maximum-weight matching problem in a graph. For the k -class classification zero-one error $\ell_{0-1}(y, y') = \mathbf{1}\{y \neq y'\}$, the loss (16.4.3) becomes

$$\phi(s, y) = \max_{j \neq y} [1 + s_j - s_y]_+.$$

JCD Comment: Make an exercise to show how the loss (16.4.3) is a maximum-weight matching.

Unfortunately, excepting (essentially) the case of binary classification, such margin-based classifiers are necessarily inconsistent. We can make this clearest in the case of binary classification.

Proposition 16.4.6. *Consider k -class binary classification with zero-one loss $\ell_{0-1}(y, y') = \mathbf{1}\{y \neq y'\}$ and let $y^* = \operatorname{argmax}_{y \in [k]} P(Y = y)$. Then for the hinge loss (16.4.3), the following hold:*

- (i) *If $P(Y = y^*) > \frac{1}{2}$, then $s = e_{y^*}$ minimizes $\mathbb{E}_P[\phi(s, Y)]$.*
- (ii) *If $P(Y = y^*) < \frac{1}{2}$, then $s = \mathbf{0}$ minimizes $\mathbb{E}_P[\phi(s, Y)]$.*

The proposition shows that surrogate consistency can occur *only* if the problem is low noise enough that the correct label has at least probability $\frac{1}{2}$. Whether this is a reasonable assumption depends on the application, as in some problems the correct label is relatively obvious, meaning we are in case (i) above. The result also helps to explain the consistency results in the binary setting, where of course the correct label necessarily has probability at least $\frac{1}{2}$.

Proof Let $p_j = P(Y = j)$ for shorthand, and define the function $g(s, y) := [s_y - s_j]_{j \neq y} \in \mathbb{R}^{k-1}$. Then

$$\mathbb{E}_P[\phi(s, Y)] = \sum_{y=1}^k p_y \left[1 - \min_j g_j(s, y) \right]_+$$

Letting $g_y(s) = \min_{j \neq y} \{s_y - s_j\}$, the objective becomes

$$\mathbb{E}_P[\phi(s, Y)] = \sum_{y=1}^k p_y [1 - g_y(s)]_+.$$

Let s^* minimize the objective and $g_y^* = g_y(s^*)$ for shorthand. We claim two properties: (P1) that $g_y^* \leq 1$ for each y and (P2) that if $\hat{y} = \operatorname{argmax}_j g_j^*$, then for $j \neq \hat{y}$, $g_j^* = -g_{\hat{y}}^*$. Assuming property (P1), the objective becomes $\sum_y p_y(1 - g_y)$, so that the problem is equivalent to maximizing $\sum_y p_y g_y$. When property (P2) holds, $\sum_y p_y g_y = (2p_{\hat{y}} - 1)g_{\hat{y}}$, this in turn is equivalent to solving

$$\underset{s}{\text{maximize}} \quad (2p_{\hat{y}} - 1)g_{\hat{y}}(s) \quad \text{subject to} \quad g_{\hat{y}} \in [0, 1].$$

Finally, we see that if no y^* satisfies $P(Y = y^*) \geq \frac{1}{2}$, then $(2p_y - 1) < 0$ for all y , and so choosing $s = \mathbf{0}$ evidently maximizes the preceding display. If $P(Y = y^*) > \frac{1}{2}$, then the choice $s = e_{y^*}$ evidently achieves the maximum above.

We return to properties (P1) and (P2). For (P1), suppose to the contrary that $g_y^* > 1$ for some y , then we have $[1 - g_y^*]_+ = 0$, while for $j \neq y$,

$$g_j^* \leq s_j^* - s_y^* = -(s_y^* - s_j^*) < -\min_j (s_y^* - s_j^*) = -g_y^* < -1,$$

which would yield objective $\mathbb{E}_P[\phi(s^*, Y)] = \sum_{j=1}^k p_j [1 - g_j^*]_+$. If $p_y < \frac{1}{2}$, this quantity satisfies $\sum_{j=1}^k p_j [1 - g_j^*]_+ > 2(1 - p_y) > 1$, so that $s = \mathbf{0}$ would give better objective; if $p_y \geq \frac{1}{2}$, then setting $s = e_y$ would yield $\mathbb{E}[\phi(s, Y)] = 2(1 - p_y) < \sum_{j \neq y} p_j [1 - g_j^*]_+$. To see property (P2), observe that

$g_{\hat{y}}^* \in [0, 1]$. We know that $g_j^* \leq s_j^* - s_{\hat{y}}^* \leq -g_{\hat{y}}^*$ as above. Consider the scores vector with $s_j = 0$ for $j \neq \hat{y}$ and $s_{\hat{y}} = g_{\hat{y}}^*$. Then $g_j(s) = -g_{\hat{y}}^* \geq g_j^*$ while $g_{\hat{y}}(s) = g_{\hat{y}}^*$. So

$$\sum_{j=1}^k p_j [1 - g_j^*]_+ = p_{\hat{y}} [1 - g_{\hat{y}}^*]_+ + \sum_{j \neq \hat{y}} p_j [1 - g_j^*]_+ \geq p_{\hat{y}} [1 - g_{\hat{y}}(s)]_+ + \sum_{j \neq \hat{y}} p_j [1 - g_j(s)]_+,$$

and the inequality is strict if $g_j^* < -g_{\hat{y}}^*$ for any j . \square

16.4.2 Structured prediction losses via the generalized entropy

Given the failure that Proposition 16.4.6 highlights for the margin-based losses (16.4.3), surrogate risk consistency in general cases will require more. To that end, we consider generalizing the entropy-based surrogates (16.3.1). Let us consider the construction of the generalized entropy associated to the loss (16.4.1), where for simplicity we (w.l.o.g.) assume $c = 0$. Then for a distribution P on $Y \in \mathcal{Y}$, we define the entropy

$$H_\ell(P) := \min_{y \in \mathcal{Y}} \mathbb{E}_P [\tau(y)^\top A \tau(Y)].$$

Let $\Delta_{\mathcal{Y}}$ be the set of probability vectors (or p.m.f.s) on \mathcal{Y} , where $p \in \Delta_{\mathcal{Y}}$ assigns probability p_y to $y \in \mathcal{Y}$. Recalling Section 14.2.3, where we considered vector-valued targets y , we generalize the mean mapping by defining *marginal polytope* associated with the part-mapping τ of y by

$$\mathcal{M} := \text{Conv}(\{\tau(y)\}_{y \in \mathcal{Y}}) = \left\{ \sum_{y \in \mathcal{Y}} p_y \tau(y) \mid p \in \Delta_{\mathcal{Y}} \right\}$$

and the mean mapping $\mu : \Delta_{\mathcal{Y}} \rightarrow \mathcal{M}$ by

$$\mu(P) := \mathbb{E}_P[\tau(Y)] = \sum_{y \in \mathcal{Y}} p_y \tau(y)$$

when P has p.m.f. p . Then we immediately see that

$$H_\ell(P) = \min_{y \in \mathcal{Y}} \tau(y)^\top A \mu(P) = \inf_{\nu \in \mathcal{M}} \nu^\top A \mu(P),$$

so that H_ℓ is equivalent for all distributions P with identical mean $\mu(P) = \mathbb{E}_P[\tau(Y)]$. With some abuse of notation, we can therefore define the negative entropy mapping by

$$\Omega(\mu) := - \min_{y \in \mathcal{Y}} \tau(y)^\top A \mu + \mathbf{I}_{\mathcal{M}}(\mu),$$

where $\mathbf{I}_{\mathcal{M}}(\mu) = 0$ if $\mu \in \mathcal{M}$ and $+\infty$ otherwise, which evidently satisfies $\Omega(\nu) = H_\ell(P)$ for any $\nu \in \mathcal{M}$ satisfying $\nu = \mu(P)$; Ω is convex.

With this construction in hand, the immediate generalization of the surrogate (16.3.1) is

$$\varphi(s, y) := -s^\top \tau(y) + \Omega^*(s), \tag{16.4.4}$$

which is closed and convex. For this surrogate, we immediately observe that

$$\mathbb{E}_p[\varphi(s, Y)] = -s^\top \mu(p) + \Omega^*(s),$$

meaning we immediately obtain the following analogue of Proposition 16.3.1:

$$\inf_s \mathbb{E}_p[\varphi(s, Y)] = \inf_s \left\{ -s^\top \mu(p) + \Omega^*(s) \right\} = -\sup_s \left\{ s^\top \mu(p) - \Omega^*(s) \right\} = -\Omega(\mu(p)) = H_\ell(p).$$

Example 16.4.7 (Multiclass classification): For k -class classification with the zero-one loss $\ell_{0-1}(y, y') = \mathbf{1}\{y \neq y'\}$, we have $\tau(y) = e_y$ and take $A = \mathbf{1}\mathbf{1}^\top - I_k$ in the representation (16.4.1). Then for $p \in \Delta_k$,

$$H_\ell(p) = \inf_{q \in \Delta_k} q^\top (\mathbf{1}\mathbf{1}^\top - I_k)p = 1 - \sup_{q \in \Delta_k} q^\top p = 1 - \max_j p_j,$$

which of course we could have obtained by direct calculation. Exercise 16.3 asks you to show that for $\Omega(p) = \max_j p_j - 1$,

$$\Omega^*(s) = 1 + \max \left\{ s_{(1)} - 1, \frac{s_{(1)} + s_{(2)} - 1}{2}, \dots, \frac{s_{(1)} + \dots + s_{(k)} - 1}{k} \right\}, \quad (16.4.5)$$

where $s_{(1)} \geq s_{(2)} \geq \dots \geq s_{(k)}$ denotes the vector s sorted in decreasing order. Then the convex surrogate loss of the zero-one loss is

$$\varphi(s, y) = 1 - s_y + \max \left\{ s_{(1)} - 1, \frac{s_{(1)} + s_{(2)} - 1}{2}, \dots, \frac{s_{(1)} + \dots + s_{(k)} - 1}{k} \right\}.$$

In passing, we note that other losses *also* generate the 0-1 entropy $H(p) = 1 - \max_j p_j$. For example, if we define

$$\varphi^{\text{hinge}}(s, y) := \sum_{j \neq y} [1 + s_j - s_y]_+,$$

then $\inf_s \mathbb{E}_p[\varphi^{\text{hinge}}(s, Y)] = k(1 - \max_j p_j)$. (Exercise 16.3 asks you to prove this as well.) \diamond

The first question is why, for the particular discrete-type loss (16.4.1), the surrogate (16.4.4) should be (surrogate-risk) consistent. To gain some intuition for this case, we present a few calculations that rely on convex analysis, before we move to the more sophisticated argument to come, which implies surrogate risk consistency. As always, we consider only pointwise versions of the risk—as surrogate risk consistency requires only this—and fix a $P \in \Delta_{\mathcal{Y}}$ and its induced $\mu = \mu(P)$. Recall the conditional surrogate loss $r_\varphi(s, P) = \mathbb{E}_P[\varphi(s, Y)]$, and consider the vector s minimizing it, thus satisfying

$$r_\varphi(s, P) = -s^\top \mu + \Omega^*(s).$$

As s minimizes r , we have $\Omega(\mu) - s^\top \mu + \Omega^*(s) = 0$, and thus (using Proposition C.2.3 in Appendix C.2.1) we necessarily have

$$s \in \partial\Omega(\mu).$$

As $\Omega(\mu) = \max_{y \in \mathcal{Y}} -\tau(y)^\top A\mu + \mathbf{I}_{\mathcal{M}}(\mu)$, then recalling the definition $\mathcal{N}_{\mathcal{M}}(\mu) = \{v \mid \langle v, \nu - \mu \rangle \leq 0 \text{ for all } \nu \in \mathcal{M}\}$ of the normal cone to \mathcal{M} at μ (see Definition C.1), we have

$$\begin{aligned} \partial\Omega(\mu) &= \text{Conv}\{-A^\top \tau(y) \mid \tau(y)^\top A\mu = \min_{y'} \tau(y')^\top A\mu\} + \mathcal{N}_{\mathcal{M}}(\mu) \\ &= \{-A^\top \nu \mid \nu^\top A\mu = H_\ell(P), \nu \in \mathcal{M}\} + \mathcal{N}_{\mathcal{M}}(\mu), \end{aligned} \quad (16.4.6)$$

by Proposition C.1.4, where we recall that $\mu = \mu(P)$. If we make the (unrealistically) simplifying assumption that μ is interior to \mathcal{M} (say, for example, if P assigns positive probability to all $y \in \mathcal{Y}$), i.e. $\mu \in \text{int } \mathcal{M}$, then $\mathcal{N}_{\mathcal{M}}(\mu) = \{\mathbf{0}\}$. If we also assume there is only a single label $y^* \in \mathcal{Y}$ minimizing $\tau(y^*)^\top A \mu$, that is, a single best prediction for the probabilities P on Y , then we obtain that $s = -A^\top \tau(y^*)$. Then of course the prediction function becomes

$$\text{pred}(s) = \operatorname{argmax}_{y \in \mathcal{Y}} \left\{ -\tau(y)^\top A^\top \tau(y^*) \right\} = y^*$$

by the assumed identifiability condition (16.4.2) on ℓ .

In summary, we have demonstrated the following proposition:

Proposition 16.4.8. *Let P be a distribution on (X, Y) such that with probability 1 over the draw of X , $\mathbb{E}[\tau(Y) \mid X] \in \text{int } \mathcal{M}$ and $\hat{y}(X) = \operatorname{argmin}_{y \in \mathcal{Y}} \mathbb{E}[\ell(y, Y) \mid X]$ is unique. Then the surrogate minimizer*

$$f^*(x) := \operatorname{argmin}_s \mathbb{E}[\varphi(s, Y) \mid X = x]$$

satisfies

$$\mathbb{E}[\ell(\text{pred}(f^*(X)), Y)] = \inf_f \mathbb{E}[\ell(f(X), Y)].$$

Thus, under appropriate regularity conditions on the distribution of the pairs (X, Y) , the surrogate loss we have constructed from the entropy is consistent, in that given infinite data, it provides optimal predictions.

To show the more sophisticated calibration guarantee of Definition 16.4, which (via Theorem 16.2.4) implies surrogate risk consistency, we will use an additional assumption.

JCD Comment: Could we do it with an alternative assumption to A.16.1 that there exists a distribution P^* for which $\mu(P^*) = \mu(P)$ and $P^*(Y = y) > 0$ for each $y \in y^*(\mu)$? Then if $\nu \in \text{Conv}\{\tau(y)\}$ in Eq. (16.4.9), we have $\nu = \mu(P')$ for some P' supported only on $y^*(\mu)$. Now there *exists* a distribution P_ν with $\mu(P_\nu) = \nu$ and $P_\nu(y) > 0$ for each $y \in y^*(\nu)$, and $y^*(\mu(P_\nu)) = y^*(\nu)$ clearly.

Assumption A.16.1. *The loss ℓ is symmetric, and if y minimizes $\mathbb{E}_P[\ell(y, Y)] = \mathbb{E}_P[\tau(y)^\top A \tau(Y)]$ then $P(Y = y) > 0$.*

The assumption is somewhat restrictive, though we can relax it at the expense of a much more sophisticated analysis. For k -class multiclass classification problems with the zero-one loss as in Example 16.4.7, we see immediately that if $y \in \operatorname{argmin}_y \mathbb{E}_P[\ell_{0-1}(y, Y)]$, then $P(Y = y) \geq \frac{1}{k} > 0$. In bipartite matching problems (Example 16.4.4), however, the assumption can fail when we allow arbitrary distributions on the collections of matchings (see Exercise 16.4).

Nonetheless, under the assumption, we have the following theorem, whose proof we provide in Section 16.4.3.

Theorem 16.4.9. *Let Assumption A.16.1 hold for a structured prediction loss ℓ taking the form (16.4.1). Then the surrogate $\varphi(s, y) = -s^\top \tau(y) + \Omega^*(s)$ is surrogate-risk consistent for ℓ .*

So the generalized entropy provides a general purpose construction that guarantees consistency. As we note above, in some cases Assumption A.16.1 may fail to hold for an arbitrary distribution P

on the set \mathcal{Y} . If we assume the problem is such that the conditional distribution of Y satisfies the low-noise condition

$$\hat{y} \in \operatorname{argmin}_{y \in \mathcal{Y}} \mathbb{E}[\ell(y, Y) \mid X = x] \text{ implies } P(Y = \hat{y} \mid X = x) > 0 \quad (16.4.7)$$

for all x except a null set, then the theorem still holds.

Corollary 16.4.10. *Let the distribution on (X, Y) satisfy the low noise condition (16.4.7). Then the surrogate $\varphi(s, y)$ is surrogate-risk consistent for ℓ , that is, $R_\varphi(f_n) \rightarrow R_\varphi^*$ implies $R(f_n) \rightarrow R^*$.*

16.4.3 Proof of Theorem 16.4.9

The proof proceeds similarly to the guarantee that $\operatorname{pred}(s)$ is correct in this case. Recall the gap functional (16.2.1),

$$\Delta_\varphi(\epsilon, P) := \inf_s \{r_\varphi(s, P) - r_\varphi^*(P) \mid r(s, P) - r^*(P) \geq \epsilon\},$$

where $r(s, P) = \mathbb{E}_P[\ell(\operatorname{pred}(s), Y)]$ and $r^*(s, P) = \inf_s \mathbb{E}_P[\ell(\operatorname{pred}(s), Y)]$. In this case, we may simplify the quantities by writing out the entropy functionals explicitly as

$$r(s, P) = \tau(\hat{y}(s))^\top A\mu - \inf_{\nu \in \mathcal{M}} \nu^\top A\mu,$$

where $\hat{y}(s)$ is (an arbitrary) element of the prediction set $\operatorname{argmax}_y s^\top \tau(y)$. We need only show that $\Delta_\varphi(\epsilon, P) > 0$ whenever $\epsilon > 0$, which we prove by contradiction.

Thus, assume for the sake of contradiction that $\Delta_\varphi(\epsilon, P) = 0$. As the losses φ are piecewise linear and the set of s such that $r(s, P) - r^*(P) \geq \epsilon$ is a union of polyhedra, there must be s achieving the infimum, and so for some vector of scores s , we have

$$\Omega^*(s) - s^\top \mu + \Omega(\mu) = 0$$

while $\hat{y}(s)$ is incorrect. Following the calculation (16.4.6), we thus obtain that for some $\nu^* \in \mathcal{M}$ satisfying $\langle \nu^*, A\mu \rangle = \min_y \tau(y)^\top A\mu$ and a vector $w \in \mathcal{N}_\mathcal{M}(\mu)$, we have

$$s = -A^\top \nu^* + w.$$

For $\nu \in \mathcal{M}$, define the shorthand $y^*(\nu) = \operatorname{argmin}_y \tau(y)^\top A\nu$, which is a set-valued mapping. If we can show the inclusions

$$\hat{y}(s) \subset y^*(\nu^*) \subset y^*(\mu(P)), \quad (16.4.8)$$

then the proof is complete, as we would evidently have our desired contradiction: necessarily, for any $\hat{y} \in \hat{y}(s)$, we would obtain $\tau(\hat{y})^\top A\mu(P) = \min_y \tau(y)^\top A\mu(P)$.

To see the inclusion $y^*(\nu^*) \subset y^*(\mu(P))$ is relatively straightforward. Let

$$\nu' \in \operatorname{Conv} \left\{ \tau(y) \mid \tau(y)^\top A\mu(P) = \min_{y' \in \mathcal{Y}} \mathbb{E}_P[\ell(y', Y)] = H_\ell(P) \right\} \quad (16.4.9)$$

be otherwise arbitrary. For all P' satisfying the mean-mapping equality $\nu' = \mu(P')$, the identifiability assumption A.16.1 guarantees that if $y \in y^*(\nu')$, we must have $P'(Y = y) > 0$. That is, $y^*(\nu')$ is contained in the supports

$$y^*(\nu') \subset \bigcap_{P'} \{\operatorname{supp} P' \mid \nu' = \mu(P')\}.$$

In equation (16.4.9), we may represent each ν' via P' supported only on $\{y \mid \tau(y)^\top A\mu(P) = H_\ell(P)\} = y^*(\mu(P))$. Thus $y^*(\nu') \subset y^*(\mu(P))$ for all ν' in the convex hull (16.4.9), and in particular for ν^* .

The first inclusion in the chain (16.4.8) is more challenging. We begin a convex analytic result that allows us to simplify maximizers of $s^\top \tau(y)$ in y .

Lemma 16.4.11. *Let $w \in \mathcal{N}_M(\mu)$ be the element satisfying $s = -A^\top \nu^* + w$. Then for any $y \in \mathcal{Y}$ and any $z \in \text{supp } P$,*

$$\langle \tau(z) - \tau(y), w \rangle \geq 0.$$

Proof Fix any $y \in \mathcal{Y}$ and let $z \in \text{supp } P$, so that $p_z = P(Y = z) > 0$. Then for some vector $\alpha \in \text{Conv}(\tau(y') \mid y \notin \{y, z\})$, we can write $\mu(P) = \lambda_y \tau(y) + \lambda_z \tau(z) + (1 - \lambda_y - \lambda_z)\alpha$, where $\lambda_y \geq 0, \lambda_z \geq p_z > 0$, and $\lambda_y + \lambda_z \leq 1$. The vector $\nu = (\lambda_y + \lambda_z)\tau(y) + (1 - \lambda_y - \lambda_z)\alpha$ similarly satisfies $\nu \in \mathcal{M}$. By the definition of the normal cone $\mathcal{N}_M(\mu)$, we know that $w^\top (\mu' - \mu) \leq 0$ for all $\mu' \in \mathcal{M}$, and in particular this holds for $\mu' = \nu$. As

$$\nu - \mu = \lambda_z(\tau(y) - \tau(z)),$$

we obtain

$$\lambda_z w^\top (\tau(y) - \tau(z)) \leq 0,$$

and as $\lambda_z > 0$ the lemma follows. \square

With Lemma 16.4.11 in hand, we can consider the prediction set $\hat{y}(s) = \text{argmax}_y s^\top \tau(y)$. As $s = -A^\top \nu^* + w$, we have

$$\hat{y}(s) = \text{argmax}_{y \in \mathcal{Y}} \left\{ -\tau(y)^\top A^\top \nu^* + \tau(y)^\top w \right\} = \text{argmax}_{y \in \mathcal{Y}} \left\{ -\tau(y)^\top A \nu^* + \tau(y)^\top w \right\},$$

where we have used the assumed symmetry of A . Let $y \in y^*(\nu^*)$ and $y' \notin y^*(\nu^*)$, so that $\tau(y')^\top A \nu^* > \tau(y)^\top A \nu^*$. Then by our earlier argument that $y^*(\nu^*) \subset y^*(\mu(P))$ and so $P(Y = y) > 0$, we obtain from Lemma 16.4.11 that

$$\tau(y')^\top w \leq \tau(y)^\top w.$$

We then see that

$$-\tau(y')^\top A \nu^* + \tau(y')^\top w < -\tau(y)^\top A \nu^* + \tau(y)^\top w,$$

and so $y' \notin \hat{y}(s)$. In particular, $\hat{y}(s) \subset y^*(\nu^*)$ as desired.

16.5 Universal loss equivalence and entropies

When we solve a statistical learning problem, we do more than simply learning a (measurable) predictor $f : \mathcal{X} \rightarrow \mathbb{R}^k$ —we frequently must estimate or learn a representation as well. Our surrogate risk consistency results typically require that we minimize over all measurable functions f , enabling us to focus on pointwise calculations conditional on X , meaning that they say little about this more elaborated process of jointly estimating a data representation and predictor. By extending our theory of losses and associated entropies, however, we can show that losses can exhibit stronger equivalences than the consistency results so far.

We re-adopt the notation of the generalized entropy associated with the loss $\ell : \mathbb{R}^k \times \mathcal{Y} \rightarrow \mathbb{R}$ we use in Chapter 14.1.3, defining

$$H_\ell(Y) := \inf_s \mathbb{E}[\ell(s, Y)]$$

and the conditional generalized entropy

$$H_\ell(Y | X) := \mathbb{E} \left[\inf_s \mathbb{E}[\ell(s, Y) | X] \right],$$

the average expected loss observing X . Then, as in definition (14.1.7), the information that X carries about Y is the amount that X reduces the expected loss in prediction Y , as measured by the loss ℓ :

$$I_\ell(X; Y) := H_\ell(Y) - H_\ell(Y | X).$$

We have seen examples in which the entropy associated with different losses was (up to a constant multiplicative factor) identical and constructed surrogate risks with the property that the associated entropy was identical to the original loss ℓ . For example, in binary classification, letting $p = P(Y = 1)$ and $(1 - p) = P(Y = -1)$, the zero-one loss has entropy

$$H_{\ell_{0-1}}(P) = \inf_{s \in \mathbb{R}} P(sY \leq 0) = 1 - \max_{y \in \{\pm 1\}} P(Y = y) = \min\{p, 1 - p\}$$

while the hinge loss $\phi(t) = [1 - t]_+$ has entropy

$$H_\phi(P) = \inf_{s \in \mathbb{R}} \{p[1 - s]_+ + (1 - p)[1 - s]_+\} = 2 \min\{p, 1 - p\}.$$

We might ask whether two losses having identical entropies guarantees something more than the basic consistency properties we have developed.

To do this, we adopt the language of quantization and data processing. We will be somewhat abstract and say that a quantizer of $X \in \mathcal{X}$ is any mapping $\mathbf{q} : \mathcal{X} \rightarrow \mathcal{Z}$ for some space \mathcal{Z} . We will think of such quantizers as a data representation, taking raw inputs x and transforming them into $z \in \mathcal{Z}$. Then for a prediction function $f : \mathcal{Z} \rightarrow \mathbb{R}^k$, we define the quantized risk and optimal quantized risks

$$L(f | \mathbf{q}) := \mathbb{E}[\ell(f(\mathbf{q}(X)), Y)] \quad \text{and} \quad L^*(\mathbf{q}) := \inf_f L(f | \mathbf{q}).$$

Rewriting these quantities in terms of the loss-based information, we have $H_\ell(Y | \mathbf{q}(X)) = \mathbb{E}[\inf_s \mathbb{E}[\ell(s, Y) | \mathbf{q}(X)]] = L^*(\mathbf{q})$, and

$$I_\ell(\mathbf{q}(X); Y) = H_\ell(Y) - H_\ell(Y | \mathbf{q}(X)) = H_\ell(Y) - L^*(\mathbf{q}).$$

Then a quantizer \mathbf{q}_1 outperforms \mathbf{q}_2 if

$$L^*(\mathbf{q}_1) < L^*(\mathbf{q}_2) \quad \text{if and only if} \quad I_\ell(\mathbf{q}_1(X); Y) > I_\ell(\mathbf{q}_2(X); Y),$$

that is, $\mathbf{q}_1(X)$ typically carries more information about Y than $\mathbf{q}_2(X)$.

We can now ask for stronger versions of surrogate risk consistency, where in addition to consistency of φ for a given loss ℓ , we ask that using φ to choose a data representation (or quantizer) from a class \mathcal{Q} of potential quantizers should be equivalent to using the original loss ℓ . We therefore define the following equivalence:

Definition 16.7. Losses ℓ_1 and ℓ_2 are universally equivalent if for all distributions on (X, Y) and all quantizers \mathbf{q}_1 and \mathbf{q}_2 ,

$$I_{\ell_1}(\mathbf{q}_1(X); Y) \leq I_{\ell_1}(\mathbf{q}_2(X); Y) \quad \text{if and only if} \quad I_{\ell_2}(\mathbf{q}_1(X); Y) \leq I_{\ell_2}(\mathbf{q}_2(X); Y).$$

We note in passing that swapping the roles of \mathbf{q}_1 and \mathbf{q}_2 and taking contrapositives, we an equivalent formulation to Definition 16.7 is that

$$I_{\ell_1}(\mathbf{q}_1(X); Y) < I_{\ell_1}(\mathbf{q}_2(X); Y) \quad \text{if and only if} \quad I_{\ell_2}(\mathbf{q}_1(X); Y) < I_{\ell_2}(\mathbf{q}_2(X); Y).$$

It is almost immediate that if H_{ℓ_1} and H_{ℓ_2} are generalized entropies associated with losses ℓ_1 and ℓ_2 equal to multiplicative constants, then the losses ℓ_1 and ℓ_2 are universally equivalent. The converse turns out to be true as well for multiclass classification problems, where $\ell : \mathbb{R}^k \times [k] \rightarrow \mathbb{R}$.

Theorem 16.5.1. Let the multiclass losses ℓ_1 and ℓ_2 be bounded below and H_1 and H_2 be the associated generalized entropies. Then ℓ_1 and ℓ_2 are universally equivalent if and only if there exist $a > 0$, $b \in \mathbb{R}^k$, and $c \in \mathbb{R}$ such that for all distributions on $Y \in [k]$,

$$H_1(Y) = aH_2(Y) + b^\top p + c, \tag{16.5.1}$$

where $p = [P(Y = y)]_{y=1}^k$ is the p.m.f. of Y .

One direction of the theorem, as we mention above, is easy: if the entropy equivalence (16.5.1) holds, then ℓ_1 and ℓ_2 are universally equivalent. Indeed, we see that

$$I_1(\mathbf{q}_1(X); Y) \leq I_1(\mathbf{q}_2(X); Y) \quad \text{if and only if} \quad H_1(Y | \mathbf{q}_1(X)) \geq H_1(Y | \mathbf{q}_2(X)),$$

and letting $p = [P(Y = y)]_{y=1}^k$ be the marginal probabilities of the label Y , the latter occurs if and only if

$$aH_2(Y | \mathbf{q}_1(X)) + b^\top [P(Y = y)]_{y=1}^k + c \geq aH_2(Y | \mathbf{q}_2(X)) + b^\top [P(Y = y)]_{y=1}^k + c.$$

The sufficiency of condition (16.5.1) is thus immediate; for its necessity, see Section 16.5.1.

Example 16.5.2 (Universal equivalence for 0-1 losses): Consider classification with the 0-1 loss $\ell_{0-1}(s, y) = \mathbf{1}\{s_y \leq \max_{j \neq y} s_j\}$. Then $H_{\ell_{0-1}}(Y) = 1 - \max_y P(Y = y)$. Several convex surrogates are both surrogate risk consistent and universally equivalent to the zero-one loss. Recalling Example 16.4.7,

$$\varphi(s, y) := 1 - s_y + \max \left\{ s_{(1)} - 1, \frac{s_{(1)} + s_{(2)} - 1}{2}, \dots, \frac{s_{(1)} + \dots + s_{(k)} - 1}{k} \right\}$$

is consistent and equivalent to ℓ_{0-1} , with $H_\varphi(Y) = H_{\ell_{0-1}}(Y)$. The hinge-type loss $\varphi^{\text{hinge}}(s, y) = \sum_{j \neq y} [1 + s_j - s_y]_+$ also satisfies $H_{\varphi^{\text{hinge}}}(Y) = k \cdot H_{\ell_{0-1}}(Y)$. \diamond

We complete this section with a corollary of Theorems 16.5.1 and 16.4.9, which shows that the entropy-based construction of surrogate losses is enough to achieve *both* surrogate risk consistency and universal equivalence of losses.

Corollary 16.5.3. Let ℓ be a structured prediction loss of the form (16.4.1) with associated generalized entropy $H_\ell(Y) = \inf_s \mathbb{E}[\ell(\text{pred}(s), Y)]$, and let Assumption A.16.1 hold. Then the surrogate (16.4.4), $\varphi(s, y) := -s^\top \tau(y) + \Omega^*(s)$, is both universally equivalent to and surrogate risk consistent for ℓ .

In brief, given a (discrete) loss, we can construct its entropy, develop a convex surrogate whose associated entropy is identical to the original function, and guarantee consistency of the convex surrogate.

16.5.1 Proof of Theorem 16.5.1

We give the proof of the “hard” direction of Theorem 16.5.1: that is, that if ℓ_1 and ℓ_2 are universally equivalent, then their associated entropies are necessarily linearly related (16.5.1). We prove the result in the slightly simpler case of binary classification, as the more general result introduces few new ideas except that it requires more technical care.

We thus work with margin-based losses $\phi_i : \mathbb{R} \rightarrow \mathbb{R}_+$, $i = 1, 2$, where without any loss of generality we assume $\inf_s \ell_i(s) = 0$ (as we may subtract a constant), and we have generalized (binary) entropies

$$h_i(p) := \inf_s \{p\phi_i(-s) + (1-p)\phi_i(s)\}.$$

By inspection, each h_i is a closed concave function (as it is the infimum of linear functions of p), and by symmetry they satisfy $h_i(0) = h_i(1) = 0$ and $h_i(\frac{1}{2}) = \sup_{p \in [0,1]} h_i(p)$. We show that these entropies satisfy a particular *order equivalence* property on $[0, 1]$, which will turn out to be sufficient to prove their equality.

To motivate what follows, recall that universal equivalence (Def. 16.7) must hold for *all* distributions on (X, Y) , and hence all (measurable) spaces \mathcal{X} and joint distributions on $\mathcal{X} \times \{\pm 1\}$. Thus, consider a space \mathcal{X} that we can partition into sets $\{A, A^c\}$ or $\{B, B^c\}$, where we take the conditional distributions

$$Y \mid X \in A = \begin{cases} 1 & \text{w.p. } p_a \\ -1 & \text{w.p. } 1 - p_a, \end{cases} \quad Y \mid X \in A^c = \begin{cases} 1 & \text{w.p. } q_a \\ -1 & \text{w.p. } 1 - q_a, \end{cases}$$

and

$$Y \mid X \in B = \begin{cases} 1 & \text{w.p. } p_b \\ -1 & \text{w.p. } 1 - p_b, \end{cases} \quad Y \mid X \in B^c = \begin{cases} 1 & \text{w.p. } q_b \\ -1 & \text{w.p. } 1 - q_b, \end{cases}$$

where we require the consistency conditions that the marginals over Y remain constant, that is, if $P(A) = P(X \in A)$, we have

$$P(A)p_a + P(A^c)q_a = P(Y = 1) = P(B)p_b + P(B^c)q_b.$$

Then evidently by defining quantizers \mathbf{q}_1 and \mathbf{q}_2 such that $\mathbf{q}_1(x) = \mathbf{1}\{x \in A\}$ and $\mathbf{q}_2(x) = \mathbf{1}\{x \in B\}$, we have

$$\begin{aligned} I_{\phi_1}(\mathbf{q}_1(X); Y) &= h_1(P(Y = 1)) - P(A)h_1(p_a) - (1 - P(A))h_1(q_a), \\ I_{\phi_1}(\mathbf{q}_2(X); Y) &= h_1(P(Y = 1)) - P(B)h_1(p_b) - (1 - P(B))h_1(q_b), \end{aligned}$$

and similarly for I_{ϕ_2} . Then universal equivalence implies that

$$\begin{aligned} P(A)h_1(p_a) + (1 - P(A))h_1(q_a) &\leq P(B)h_1(p_b) + (1 - P(B))h_1(q_b) \quad \text{if and only if} \\ P(A)h_2(p_a) + (1 - P(A))h_2(q_a) &\leq P(B)h_2(p_b) + (1 - P(B))h_2(q_b) \end{aligned}$$

whenever the consistency condition $P(A)p_a + (1 - P(A))q_a = P(B)p_b + (1 - P(B))q_b$ holds. As we may choose \mathcal{X} and the probabilities, we can take $P(A) = P(B) = \frac{1}{2}$ (so that their are two equiprobable partitions), and the preceding conditions become

$$h_1(p_a) + h_1(q_a) \leq h_1(p_b) + h_1(q_b) \quad \text{if and only if} \quad h_2(p_a) + h_2(q_a) \leq h_2(p_b) + h_2(q_b)$$

whenever $p_a + q_a = p_b + q_b$.

Generalizing this construction by taking distributions over \mathcal{X} that partition it into k equiprobable sets $\{A_1, \dots, A_k\}$ or $\{B_1, \dots, B_k\}$, each with $P(A_i) = P(B_i) = 1/k$, we see that universal equivalence implies that for any vectors $p \in [0, 1]^k$ and $q \in [0, 1]^k$ satisfies $\mathbf{1}^\top p = \mathbf{1}^\top q$,

$$\sum_{i=1}^k h_1(p_i) \leq \sum_{i=1}^k h_1(q_i) \quad \text{if and only if} \quad \sum_{i=1}^k h_2(p_i) \leq \sum_{i=1}^k h_2(q_i). \quad (16.5.2)$$

We shall give condition (16.5.2) a name, as it implies certain equivalence properties for convex functions (we can replace h_i with $-h_i$ and obtain convex functions).

Definition 16.8. Let $\Omega \subset \mathbb{R}$ be a closed interval and let $f_1, f_2 : \Omega \rightarrow \mathbb{R}$ be closed convex functions. Then f_1 and f_2 are order equivalent if for any $k \in \mathbb{N}$ and vectors $s \in \Omega^k$ and $t \in \Omega^k$ satisfying $\mathbf{1}^\top s = \mathbf{1}^\top t$, we have

$$\sum_{i=1}^k f_1(s_i) \leq \sum_{i=1}^k f_1(t_i) \quad \text{if and only if} \quad \sum_{i=1}^k f_2(s_i) \leq \sum_{i=1}^k f_2(t_i)$$

As in the brief remark following Definition 16.7, by taking complements we have as well that

$$\sum_{i=1}^k f_1(s_i) < \sum_{i=1}^k f_1(t_i) \quad \text{if and only if} \quad \sum_{i=1}^k f_2(s_i) < \sum_{i=1}^k f_2(t_i)$$

The theorem will then be proved if we can show the following lemma.

Lemma 16.5.4. Let f_1 and f_2 be order equivalent on Ω . Then there exist $a > 0$, and $b, c \in \mathbb{R}$ such that $f_1(t) = af_2(t) + bt + c$ for all $t \in \Omega$.

The proof of Lemma 16.5.4 is somewhat involved, and we proceed in three parts. The key is that order equivalence actually implies a strong relationship between *affine* combinations of points in the domain of the functions f_i , not just convex combinations of points, which guarantees that we can predict values of $f_2(v)$ for any $v \in \Omega$ by just three values of f_i evaluate in Ω . We state this as a lemma, whose proof we defer temporarily to Sec. 16.5.2

Lemma 16.5.5. If f_1 and f_2 are order equivalent on Ω , then for any $\lambda \in \mathbb{R}^k$ satisfying $\lambda^\top \mathbf{1} = 1$ and any $u \in \Omega^k$, if $\lambda^\top u = v \in \Omega$ then

$$\sum_{i=1}^k \lambda_i f_1(u_i) \leq f_1(v) \quad \text{if and only if} \quad \sum_{i=1}^k \lambda_i f_2(u_i) \leq f_2(v),$$

and the statement still holds with both inequalities replaced with strict inequalities.

In particular, if

$$\sum_{i=1}^k \lambda_i f_1(u_i) = f_1(v) \quad \text{then necessarily} \quad \sum_{i=1}^k \lambda_i f_2(u_i) = f_2(v) \quad (16.5.3)$$

whenever $\lambda \in \mathbb{R}^k$ satisfies $\lambda^\top \mathbf{1} = 1$ and $u^\top \lambda = \sum_{i=1}^k \lambda_i u_i = v$.

Second, we recognize that we may assume both f_1 and f_2 are nonlinear in the lemma; otherwise, it is immediate. Nonlinearity guarantees that

Lemma 16.5.6. *Let f be convex on \mathbb{R} . Let $u_0 < u_1$ and for $\lambda \in [0, 1]$, define $u_\lambda = (1 - \lambda)u_0 + \lambda u_1$. If there exists any $\lambda \in (0, 1)$ such that $f(u_\lambda) = \lambda f(u_0) + (1 - \lambda)f(u_1)$, then f is linear on $[u_0, u_1]$.*

We leave the proof (an algebraic manipulation using the definitions of convexity) as Question 16.11.

The last intermediate step we require in the proof of Lemma 16.5.4 is that at three particular points in the domain Ω , we can satisfy Lemma 16.5.4.

Lemma 16.5.7. *Let f_1, f_2 be order equivalent on $\Omega = [u_0, u_1]$ and $u_c = \frac{1}{2}(u_0 + u_1)$. There are $a > 0$ and $b, c \in \mathbb{R}$ such that $f_1(t) = af_2(t) + bt + c$ for $t \in \{u_0, u_c, u_1\}$.*

We can now finalize the proof of Lemma 16.5.4:

Proof Without loss of generality by an affine rescaling, we can assume that $f_1(t) = f_2(t)$ for $t \in \{u_0, u_c, u_1\}$, and our goal will be to show that $f_1(t) = f_2(t)$ for all $t \in \Omega$.

Let $v \in \Omega$ with $v \notin \{u_0, u_c, u_1\}$, and $u = [u_0 \ u_c \ u_1]^\top$ for shorthand. We seek $\lambda = (\lambda_0, \lambda_c, \lambda_1) \in \mathbb{R}^3$, where $\lambda^\top \mathbf{1} = 1$, such that both $\lambda^\top u = v$ and $\lambda_0 f_1(u_0) + \lambda_c f_1(u_c) + \lambda_1 f_1(u_1) = f_1(v)$. If we can find such a λ , then equality (16.5.3) guarantees that $f_1(v) = f_2(v)$, and we are done. As the points $(u_i, f_1(u_i))_{i=1}^3$ are not collinear (recall Lemma 16.5.6), the matrix

$$A := \begin{bmatrix} 1 & 1 & 1 \\ u_0 & u_c & u_1 \\ f_1(u_0) & f_1(u_c) & f_1(u_1) \end{bmatrix}$$

is full rank. In particular, there is a vector λ solving

$$A\lambda = [1 \ v \ f_1(v)]^\top, \quad \text{i.e. } \lambda = A^{-1} [1 \ v \ f_1(v)]^\top,$$

which evidently satisfies our desiderata. Thus $f_1(v) = f_2(v)$, and as v was arbitrary, the proof is complete. \square

16.5.2 Proof of Lemma 16.5.5

We prove the result first for λ with rational entries, as a continuity argument will give the rest. For each i , let $\alpha_i = [\lambda_i]_+$ and $\beta_i = [-\lambda_i]_+$ be the positive and negative parts of λ , so that $\lambda = \alpha - \beta$. Let $k \in \mathbb{N}$ be such that we can write $\alpha_i = \frac{a_i}{k}$ and $\beta_i = \frac{b_i}{k}$, where $a_i, b_i \in \mathbb{N}$. Then we have

$$\alpha^\top u = v + \beta^\top u \quad \text{or} \quad a^\top u = kv + b^\top u,$$

where $\mathbf{1}^\top a = k + \mathbf{1}^\top b$, as $\mathbf{1}^\top \lambda = \frac{1}{k} \mathbf{1}^\top (a - b) = 1$. Then we may define the two vectors

$$s = \underbrace{[u_1 \ \cdots \ u_1]}_{a_1 \text{ times}} \cdots \underbrace{[u_m \ \cdots \ u_m]}_{a_m \text{ times}}^\top \quad \text{and} \quad t = \underbrace{[v \ \cdots \ v]}_{k \text{ times}} \underbrace{[u_1 \ \cdots \ u_1]}_{b_1 \text{ times}} \cdots \underbrace{[u_m \ \cdots \ u_m]}_{b_m \text{ times}}^\top.$$

Then each has entries in Ω , and we have $\mathbf{1}^\top t = \mathbf{1}^\top s$. Then order equivalence (Def. 16.8) guarantees that

$$\sum_{i=1}^m a_i f_1(u_i) \leq k f_1(v) + \sum_{i=1}^m b_i f_1(u_i) \quad \text{if and only if} \quad \sum_{i=1}^m a_i f_2(u_i) \leq k f_2(v) + \sum_{i=1}^m b_i f_2(u_i)$$

and (as per the remark following the definition)

$$\sum_{i=1}^m a_i f_1(u_i) = k f_1(v) + \sum_{i=1}^m b_i f_1(u_i) \text{ if and only if } \sum_{i=1}^m a_i f_2(u_i) = k f_2(v) + \sum_{i=1}^m b_i f_2(u_i).$$

These two displays are equivalent to $\sum_{i=1}^m \lambda_i f_j(u_i) \leq f_j(v)$ and $\sum_{i=1}^m \lambda_i f_j(u_i) = f_j(v)$, respectively, for $j = 1, 2$.

We have therefore proved Lemma 16.5.5 for λ taking rational values. Because closed convex functions on \mathbb{R} are continuous on their domains, the result extends to real-valued λ .

16.5.3 Proof of Lemma 16.5.7

If either of f_1 or f_2 is linear, the other is as well, and the proof becomes trivial, so we assume w.l.o.g. they are both nonlinear.

Without loss of generality, we take $u_0 = 0$, $u_1 = 1$, and $u_c = \frac{1}{2}$ by scaling. Then we must solve the three equations

$$f_1(0) = a f_2(0) + c, \quad f_1\left(\frac{1}{2}\right) = a f_2\left(\frac{1}{2}\right) + \frac{b}{2} + c, \quad f_1(1) = a f_2(1) + b + c.$$

From the first we obtain $c = f_1(0) - a f_2(0)$, and substituting this into the third yields $b = f_1(1) - f_1(0) - a(f_2(1) - f_2(0))$. Finally, substituting both equalities into the equality with $f_1(\frac{1}{2})$ yields that

$$\begin{aligned} f_1\left(\frac{1}{2}\right) &= a \left[f_2\left(\frac{1}{2}\right) - \frac{f_2(1) - f_2(0)}{2} \right] + \frac{f_1(1) - f_1(0)}{2} + f_1(0) - a f_2(0) \\ &= a \left[f_2\left(\frac{1}{2}\right) - \frac{f_2(1) + f_2(0)}{2} \right] + \frac{f_1(1) + f_1(0)}{2}. \end{aligned}$$

As we know that f_1, f_2 are nonlinear, Lemma 16.5.6 applies, so that the convexity gaps $f_1(\frac{1}{2}) - \frac{f_1(1)+f_1(0)}{2}$ and $f_2(\frac{1}{2}) - \frac{f_2(1)+f_2(0)}{2}$ are both positive, and thus we take

$$a = \frac{f_1(\frac{1}{2}) - \frac{f_1(0)+f_1(1)}{2}}{f_2(\frac{1}{2}) - \frac{f_2(0)+f_2(1)}{2}} > 0.$$

16.6 Bibliography

Section 16.1 is based on Bartlett, Jordan, and McAuliffe [19]. See also Zhang [196], Lugosi and Vayatis [142], Zhang and Yu [197]. The general treatment follows [173]. The particular multiclass results we present are due to Zhang [195].

Point to full proof of Theorem 16.5.1.

Cite Nowak-Vila et al. [151] and for inconsistency (Prop. 16.4.6) cite Liu [140].

16.7 Exercises

Exercise 16.1: Find the suboptimality function Δ_ϕ and ψ -transform (16.1.2) for the binary classification problem ($Y \in \{-1, 1\}$) with the following losses.

- (a) Logistic loss $\phi(s) = \log(1 + e^{-s})$.
- (b) The squared hinge loss $\phi(s) = [1 - s]_+^2$.
- (c) Squared error (ordinary regression). The surrogate loss in this case for the pair (x, y) is $\frac{1}{2}(f(x) - y)^2$. Show that for $y \in \{-1, 1\}$, this can be written as a margin-based loss, and compute the associated suboptimality function Δ_ϕ and ψ -transform. Is the squared error classification calibrated?

Exercise 16.2: Suppose we have a regression problem with data (independent variables) $x \in \mathcal{X}$ and $y \in \mathbb{R}$. We wish to find a predictor $f : \mathcal{X} \rightarrow \mathbb{R}$ minimizing the probability of being far away from the true y , that is, for some $c > 0$, our loss is of the form

$$\ell(f(x), y) = \mathbf{1}\{|y - f(x)| \geq c\}.$$

Show that no loss of the form $\varphi(s, y) = |s - y|^p$, where $p \geq 1$, is Fisher consistent for the loss ℓ , even if the distribution of Y conditioned on $X = x$ is symmetric about its mean $\mathbb{E}[Y | X]$. That is, show there exists a distribution on pairs X, Y such that the set of minimizers of the surrogate

$$R_\varphi(f) := \mathbb{E}[\varphi(f(X), Y)]$$

is not included in the set of minimizers of the true risk, $R(f) = \mathbb{P}(|Y - f(X)| \geq c)$, even if the distribution of Y (conditional on X) is symmetric.

Exercise 16.3 (Hinge-type losses and entropies [76]):

- (a) Prove the equality (16.4.5) in Example 16.4.7.
- (b) Show that $\inf_s \mathbb{E}_p[\varphi^{\text{hinge}}(s, Y)] = k(1 - \max_j p_j)$.

Exercise 16.4: Let ℓ be the Hamming loss on bipartite matchings, so that for $y, y' \in \mathfrak{S}_k$, we define $\ell(y, y') = \sum_{i=1}^k \mathbf{1}\{y(i) \neq y'(i)\}$ as in Example 16.4.4. Let M_π be the permutation matrix associated with a permutation π .

- (a) Show that there is a collection of k (non-identity) permutations π_1, \dots, π_k such that

$$\frac{1}{k} \sum_{i=1}^k M_{\pi_i} = \frac{1}{k} \mathbf{1}\mathbf{1}^\top.$$

- (b) Show that Assumption A.16.1 must fail for bipartite matchings.
- (c) Using permutations that induce block diagonal matrices of the form

$$M = \begin{bmatrix} I_j & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & I_{k-j-2} \end{bmatrix},$$

show a stronger failure of Assumption A.16.1: there exist distributions on permutations for which there exists a unique minimizer of $y^* = \arg\min_y \mathbb{E}[\ell(y, Y)]$ but $P(Y = y^*) = 0$.

Exercise 16.5 (Empirics of classification calibration): In this problem you will compare the performance of hinge loss minimization and an ordinary linear regression in terms of classification performance. Specifically, we compare the performance of the hinge surrogate loss with the regression surrogate when the data is generated according to the model

$$y = \text{sign}(\langle \theta^*, x \rangle + \sigma Z), \quad Z \sim \mathcal{N}(0, 1) \quad (16.7.1)$$

where $\theta^* \in \mathbb{R}^d$ is a fixed vector, $\sigma \geq 0$ is an error magnitude, and Z is a standard normal random variable. We investigate the model (16.7.1) with a simulation study.

Specifically, we consider the following set of steps:

- (i) Generate two collections of n datapoints in d dimensions according to the model (16.7.1), where $\theta \in \mathbb{R}^d$ is chosen (ahead of time) uniformly at random from the sphere $\{\theta \in \mathbb{R}^d : \|\theta\|_2 = R\}$, and where each $x_i \in \mathbb{R}^d$ is chosen as $\mathcal{N}(0, I_{d \times d})$. Let (x_i, y_i) denote pairs from the first collection and $(x_i^{\text{test}}, y_i^{\text{test}})$ pairs from the second.

- (ii) Set

$$\hat{\theta}_{\text{hinge}} = \underset{\theta: \|\theta\|_2 \leq R}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n [1 - y_i \langle x_i, \theta \rangle]_+$$

and

$$\hat{\theta}_{\text{reg}} = \underset{\theta}{\text{argmin}} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle x_i, \theta \rangle)^2 = \underset{\theta}{\text{argmin}} \|X\theta - y\|_2^2.$$

- (iii) Evaluate the 0-1 error rate of the vectors $\hat{\theta}_{\text{hinge}}$ and $\hat{\theta}_{\text{reg}}$ on the held-out data points $\{(x_i^{\text{test}}, y_i^{\text{test}})\}_{i=1}^n$.

Perform the preceding steps (i)–(iii), using any $n \geq 100$ and $d \geq 10$ and a radius $R = 5$, for different standard deviations $\sigma = \{0, 1, \dots, 10\}$; perform the experiment a number of times. Give a plot or table exhibiting the performance of the classifiers learned on the held-out data. How do the two compare? Given that for the hinge loss we know $\Delta_\phi(\delta) = \delta$ (as presented in class), what would you expect based on the answer to Question 16.1?

I have implemented (in the `julia` language; see <http://julialang.org/>) methods for solving the hinge loss minimization problem with stochastic gradient descent so that you do not need to. The file is available at [this link](#). The code should (hopefully) be interpretable enough that if `julia` is not your language of choice, you can re-implement the method in an alternative language.

Exercise 16.6: **JCD Comment:** Check the constants and things here

Show that in the case of binary classification with the zero-one loss, so that $\varphi(s, y) = \phi(sy)$, the uniform gap (16.2.2) and binary gap (16.1.1) are equal, that is, $\Delta_\varphi(\epsilon) = \Delta_\phi(\epsilon)$. In particular, uniform calibration (Definition 16.3) and classification calibration (Definition 16.2) are equivalent.

Exercise 16.7: We generalize Theorem 16.2.4 to the case that the loss ℓ is not uniformly bounded. Assume there exists an upper bound function $B : \mathcal{X} \rightarrow \mathbb{R}_+$ such that $\mathbb{E}[B(X)] < \infty$ and $r(s, P_x) \leq r^*(P_x) + B(x)$ for all $x \in \mathcal{X}$ and $s \in \mathbb{R}^k$. Show that if φ is calibrated for the loss ℓ (Definition 16.4) if and only if it is surrogate risk consistent (Definition 16.1).

Exercise 16.8: Let $\phi(t) = [1 - t]_+$. Show by example that the surrogate (16.2.3) with $\varphi(s, y) = \sum_{j=1}^k \phi(sy - s_j)$ may be inconsistent.

Exercise 16.9: Prove Proposition 16.2.7.

JCD Comment: Outline this a bit more.

Exercise 16.10 (Bayes risk gaps): Consider a general binary classification problem with $(X, Y) \in \mathcal{X} \times \{-1, 1\}$. Let $\phi(t) = \log(1 + e^{-t})$, so that we use the logistic loss. Show that the surrogate risk gap

$$H_\phi(Y) - H_\phi(Y | X) = I(X; Y),$$

where I is the mutual information.

Exercise 16.11: Prove Lemma 16.5.6. *Hint:* without loss of generality, you may take $u_0 = 0$, $u_1 = 1$. Then for any $u \in [0, 1]$, write λ as a convex combination of either $\{0, u\}$ or $\{u, 1\}$ and use the definition of convexity.

JCD Comment: Add some exercises around f -divergences and risk gaps, which give similar flavor. Nothing too hard. Add commentary after Theorem that these exist.

Part IV

Online game playing and compression

Chapter 17

Stochastic and online convex optimization

In Chapter 5, we discussed generalization guarantees, which followed from the various uniform concentration inequalities we developed, and applications in estimation, with finite-sample guarantees for recovering parameters θ in well-behaved models and smooth loss functions. In many cases, however, we care only about minimizing some expected loss—the particular parameters may be unimportant—as when we perform large-scale machine learning or prediction tasks. In other cases, we have *online* problems, where we do not even make particular statistical assumptions, but assume simply that data arrives in a stream, and we wish to make predictions on this stream that are at least as good as some reference static predictor. Finally, the smoothness assumptions on which our convergence results of Chapter 5.3 repose—which are, in a sense, necessary—lie in opposition to some of our surrogate consistency results in Chapter 16, which suggest several optimality properties of hinge-type and other non-smooth losses, which are beyond the purview of our earlier convergence guarantees. This chapter addresses these problems.

In online optimization, we consider the following two player sequential game: we have a space Θ in which we—the learner or first player—can play points $\theta_1, \theta_2, \dots$, while nature plays a sequence of loss functions $\ell_t : \Theta \rightarrow \mathbb{R}$. The goal is to guarantee that the regret

$$\text{Reg}_n(\theta^*) := \sum_{t=1}^n [\ell_t(\theta_t) - \ell_t(\theta^*)] \quad (17.0.1)$$

grows at most sub-linearly with n for any fixed reference $\theta^* \in \Theta$. We will typically consider scenarios in which the sequence of losses ℓ_t are convex or can be appropriately modeled by convex losses, and Θ is a convex subset of \mathbb{R}^d . This defines the *online convex optimization* problem. The related problem of *stochastic convex optimization* follows when nature cannot be so capricious: instead of losses ℓ_t chosen (essentially) arbitrarily and adversarially, we have the risk minimization problems of Chapter 5. Thus, for a distribution P on $Z \in \mathcal{Z}$, we wish to minimize

$$L(\theta) := \mathbb{E}_P[\ell(\theta, Z)] \quad (17.0.2)$$

given a sample $Z_1^n \stackrel{\text{iid}}{\sim} P$, where for each $z \in \mathcal{Z}$, the loss function $\ell(\theta, z)$ is convex in θ .

The problem (17.0.2) is a main motivator of the online problem (17.0.1), as we shall see, as any algorithm that provides a strong regret guarantee for (17.0.1) implies a convergence rate for optimizing (17.0.2). Indeed, let $Z_t \stackrel{\text{iid}}{\sim} P$, and in the formulation (17.0.1), set $\ell_t(\theta) = \ell(\theta, Z_t)$. Then

by convexity, the average vector $\bar{\theta}_n := \frac{1}{n} \sum_{t=1}^n \theta_t$ satisfies

$$L(\bar{\theta}_n) \leq \frac{1}{n} \sum_{t=1}^n L(\theta_t) \stackrel{(*)}{=} \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\ell_t(\theta_t)] = \mathbb{E}\left[\frac{1}{n} \sum_{t=1}^n \ell_t(\theta_t)\right], \quad (17.0.3)$$

where the equality $(*)$ follows because θ_t is a function only of Z_1^{t-1} , and so conditionally independent of Z_t . Inequality (17.0.3) is our first example of an *online-to-batch* conversion: any convergence guarantee for the average regret in the online setting (17.0.1) implies a matching guarantee for stochastic optimization.

In the remainder of the chapter, we develop some of the modern perspective on algorithms for solving online and stochastic convex optimization problems, providing their guarantees. It will turn out that Bregman divergences, which we have developed in the context of prediction problems, provide one of our central tools. We will also develop a collection of lower bounds showing the (minimax) optimality of these methods.

17.1 Preliminaries on convex optimization

While convexity has been a touchstone throughout this text, because our focus in this chapter will be on procedures for solving convex problems and to make things self-contained, it will behoove us to review some of the relevant definitions briefly. (We refer as always to Appendix B for proofs associated to convex sets and functions.) Recall that Θ is *convex* if for all $\lambda \in [0, 1]$ and $\theta, \theta' \in \Theta$, we have

$$\lambda\theta + (1 - \lambda)\theta' \in \Theta.$$

A function f is *convex* if

$$f(\lambda\theta + (1 - \lambda)\theta') \leq \lambda f(\theta) + (1 - \lambda)f(\theta')$$

for all $\lambda \in [0, 1]$ and θ, θ' , where $f(\theta) = +\infty$ for $\theta \notin \text{dom } f$. The *subgradient set*, or *subdifferential*, of a convex function f at the point θ is

$$\partial f(\theta) := \{g \in \mathbb{R}^d \mid f(\theta') \geq f(\theta) + \langle g, \theta' - \theta \rangle \text{ for all } \theta'\},$$

and any vector $g \in \partial f(\theta)$ is a *subgradient*. For convex f , the subdifferential $\partial f(\theta)$ is non-empty for any $\theta \in \text{int dom } f$. (See Theorem B.3.3 in Appendix B.3, which shows that $\partial f(\theta)$ is non-empty on the relative interior of $\text{dom } f$.) If f is differentiable at θ , then $\partial f(\theta) = \{\nabla f(\theta)\}$ is a singleton, and we will abuse notation to simply write $\partial f(\theta) = \nabla f(\theta)$.

We now give several examples of convex functions, losses, and corresponding subgradients. The first two highlight binary classification problems, where we wish to predict associated labels $y \in \{-1, 1\}$ from a data point $x \in \mathbb{R}^d$. We focus on the online setting (17.0.1).

Example 17.1.1 (Support vector machines): We receive data in pairs $(x_t, y_t) \in \mathbb{R}^d \times \{-1, 1\}$, and the loss function

$$\ell_t(\theta) = [1 - y_t \langle \theta, x_t \rangle]_+ = \max\{1 - y_t \langle \theta, x_t \rangle, 0\},$$

which is convex because it is the maximum of two linear functions. The subgradient set is

$$\partial \ell_t(\theta) = \begin{cases} -y_t x_t & \text{if } y_t \langle \theta, x_t \rangle < 1 \\ -\lambda \cdot y_t x_t & \text{for } \lambda \in [0, 1] \text{ if } y_t \langle \theta, x_t \rangle = 1 \\ 0 & \text{otherwise} \end{cases}$$

as $s \mapsto [s]_+$ is differentiable except at $s = 0$. \diamond

Example 17.1.2 (Logistic regression): As in the support vector machine, we receive data in pairs $(x_t, y_t) \in \mathbb{R}^d \times \{-1, 1\}$, and the loss function is

$$\ell_t(\theta) = \log(1 + \exp(-y_t \langle x_t, \theta \rangle)),$$

with singleton subdifferential $\partial \ell_t(\theta) = -\frac{1}{1 + e^{y_t \langle x_t, \theta \rangle}} y_t x_t$, because the loss is \mathcal{C}^∞ . \diamond

Example 17.1.3 (Expert prediction and zero-one error): By randomization, it is possible to cast certain non-convex optimization problems as convex. Consider a problem in which there are d experts, each of which makes a prediction $x_{t,j}$ (for $j = 1, \dots, d$) at time t , represented by the vector $x_t \in \mathbb{R}^d$, of a label $y_t \in \{-1, 1\}$. Each also suffers the (non-convex) zero-one loss $\ell_{0-1}(x_{t,j}, y_t) = \mathbf{1}\{x_{t,j} y_t \leq 0\}$. Assign a weight $w_j \geq 0$ to each expert $x_{t,j}$, where the weights satisfy $\langle w, \mathbf{1} \rangle = 1$. Then if we choose a prediction $\hat{y} = \text{sign}(x_{t,j})$ with probability j , the probability of a mistake is

$$\ell_t(w) = \mathbb{P}(\hat{y} \neq y_t) = \sum_{j=1}^d w_j \ell_{0-1}(x_{t,j}, y_t) = \langle g_t, w \rangle,$$

where we have defined the vector $g_t = [\ell_{0-1}(x_{t,j}, y_t)]_{j=1}^d \in \{0, 1\}^d$. Notably, the expected zero-one loss is convex (even linear), so that its online minimization falls into the online convex programming framework. \diamond

We will see that, in spite of their frequently simplicity, online convex programming approaches have applications beyond the initial regret formulation (17.0.1).

17.2 Online convex optimization methods

The basic approach to online (and stochastic) optimization is a relatively simple one: at the t th iteration of the process, we

- i. construct a simple model (or approximation) to the instantaneous objective ℓ_t
- ii. minimize this model, regularizing so that the updated point does not move too far (or too aggressively follow spurious information)

Let us make this more precise before giving the actual algorithms and derived procedures.

Definition 17.1. A model of the loss ℓ at a point θ_0 is a function $\hat{\ell}(\cdot; \theta_0)$ satisfying

- i. **Model convexity.** The function $\theta \mapsto \hat{\ell}(\theta; \theta_0)$ is convex and subdifferentiable
- ii. **Lower bound.** The model satisfies $\hat{\ell}(\theta; \theta_0) \leq \ell(\theta)$, with equality at $\theta = \theta_0$.

The most common model is the first-order model, which takes $g \in \partial \ell(\theta_0)$ and sets

$$\hat{\ell}(\theta; \theta_0) := \ell(\theta_0) + \langle g, \theta - \theta_0 \rangle,$$

for which it is immediate that both convexity (because $\hat{\ell}$ is affine) and the lower bound condition hold. A somewhat less trivial model holds when we know that the losses are nonnegative; in this case, with $g \in \partial\ell(\theta_0)$ as before, the truncated model

$$\hat{\ell}(\theta; \theta_0) := [\ell(\theta_0) + \langle g, \theta - \theta_0 \rangle]_+$$

satisfies all the desired conditions as well.

As a brief remark, we note that the conditions in Definition 17.1 guarantee that $\hat{\ell}$ looks locally like ℓ ; in particular, we always have

$$\partial\hat{\ell}(\theta_0; \theta_0) \subset \partial\ell(\theta_0), \quad (17.2.1)$$

and in particular, if ℓ is differentiable at θ_0 then $\hat{\ell}$ is as well, with identical derivative. To see the inclusion (17.2.1), note that if $g_0 \in \partial\hat{\ell}(\theta_0; \theta_0)$, then for any θ , we have

$$\ell(\theta) \geq \hat{\ell}(\theta; \theta_0) \geq \hat{\ell}(\theta_0; \theta_0) + \langle g_0, \theta - \theta_0 \rangle = \ell(\theta_0) + \langle g_0, \theta - \theta_0 \rangle$$

by repeated application of Definition 17.1, so that $g_0 \in \partial\ell(\theta_0)$.

JCD Comment: Put a figure here about the modeling approach

17.2.1 Projected subgradient methods

Let us give a special case of the general mirror-descent-type algorithms we consider. This first is a variant of (*projected*) *online gradient descent*, and requires only that we specify a sequence $\eta_t > 0$ of non-increasing stepsizes. Here, we use the simple linear model

$$\hat{\ell}_t(\theta; \theta_t) = \ell_t(\theta_t) + \langle g_t, \theta - \theta_t \rangle,$$

and we will present a preliminary analysis because the simplicity of the update lends itself to an elementary proof of convergence.

Input: Parameter space Θ , stepsize sequence η_t .

Repeat: for each iteration t ,

i. predict $\theta_t \in \Theta$, receive function ℓ_t and suffer loss $\ell_t(\theta_t)$.

ii. let $g_t \in \partial\ell_t(\theta_t)$ and update

$$\theta_{t+1} = \operatorname{argmin}_{\theta \in \Theta} \left\{ \langle g_t, \theta \rangle + \frac{1}{2\eta_t} \|\theta - \theta_t\|_2^2 \right\}. \quad (17.2.2)$$

Figure 17.1: Projected gradient descent.

The update (17.2.2) makes clear that we trade between improving performance on ℓ_t via the linear approximation $\ell_t(\theta) \approx \ell_t(\theta_t) + \langle g_t, \theta - \theta_t \rangle$ and remaining close to θ_t according to the Euclidean distance $\|\cdot\|_2$. When we use the standard linear approximation as the model $\hat{\ell}$, that is, for some $g_t \in \partial\ell_t(\theta)$ use $\hat{\ell}_t(\theta) = \ell_t(\theta_t) + \langle g_t, \theta - \theta_t \rangle$, we may rewrite the update (17.2.2) as the two steps

$$\theta_{t+\frac{1}{2}} = \theta_t - \eta_t g_t, \quad \theta_{t+1} = \operatorname{Proj}_{\Theta}(\theta_{t+\frac{1}{2}}) = \operatorname{argmin}_{\theta \in \Theta} \left\{ \|\theta - \theta_{t+\frac{1}{2}}\|_2^2 \right\},$$

where Proj_Θ denotes Euclidean projection onto Θ . The update (17.2.2) admits an elegant analysis, which we include even though it is a special case of the more general Theorem 17.2.9 to come, because it highlights the main ideas. We use a fixed stepsize for simplicity.

Proposition 17.2.1 (Convergence of projected gradient descent). *Let $\eta_t = \eta > 0$ for all t . Then for any $\theta \in \Theta$,*

$$\sum_{t=1}^n [\ell_t(\theta_t) - \ell_t(\theta)] \leq \frac{1}{2\eta} \|\theta_1 - \theta\|_2^2 + \frac{\eta}{2} \sum_{t=1}^n \|g_t\|_2^2.$$

Proof The only important step is to write an appropriate measure of convergence: the error $\|\theta_{t+1} - \theta\|_2^2$, from which we can derive a “one-step” progress guarantee, which we then recurse more or less via a bit of algebra. To that end, note that for any $\theta \in \Theta$ and any vector v ,

$$\|\text{Proj}_\Theta(v) - \theta\|_2^2 \leq \|v - \theta\|_2^2,$$

because projections decrease (Euclidean) distance. (See Corollary B.1.13 in Appendix B.1.2.) So we have the one-step progress guarantee

$$\|\theta_{t+1} - \theta\|_2^2 \leq \|\theta_t - \eta g_t - \theta\|_2^2 = \|\theta_t - \theta\|_2^2 - 2\eta \langle g_t, \theta_t - \theta \rangle + \eta^2 \|g_t\|_2^2.$$

By the first-order conditions for convexity, we have $\ell_t(\theta_t) + \langle g_t, \theta - \theta_t \rangle \leq \ell_t(\theta)$, or $-\langle g_t, \theta_t - \theta \rangle \leq \ell_t(\theta) - \ell_t(\theta_t)$. Substituting gives

$$\|\theta_{t+1} - \theta\|_2^2 \leq \|\theta_t - \theta\|_2^2 - 2\eta [\ell_t(\theta_t) - \ell_t(\theta)] + \eta^2 \|g_t\|_2^2.$$

Rearranging and dividing by $2\eta > 0$ we obtain

$$\ell_t(\theta_t) - \ell_t(\theta) \leq \frac{1}{2\eta} \left[\|\theta_t - \theta\|_2^2 - \|\theta_{t+1} - \theta\|_2^2 \right] + \frac{\eta}{2} \|g_t\|_2^2. \quad (17.2.3)$$

Sum inequality (17.2.3) and note that the sum telescopes to obtain the proposition. \square

Typically, one imposes boundedness conditions on θ and $\|g_t\|_2$, which then allows a more evocative guarantee. As we have focused on Euclidean-type updates, we present one such prototypical result, with more discussion in Section 17.2.4.

Corollary 17.2.2. *Assume that $\Theta \subset \{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq R_2\}$ and that $\|g\|_2 \leq G_2$ for any $g \in \partial \ell_t(\theta)$ for all t . Take $\eta = \frac{R_2}{G_2} \frac{1}{\sqrt{n}}$ and let $\theta_1 = \mathbf{0}$. Then for all $\theta^* \in \Theta$,*

$$\sum_{t=1}^n [\ell_t(\theta_t) - \ell_t(\theta^*)] \leq R_2 G_2 \sqrt{n}.$$

Such guarantees—order \sqrt{n} regret, with a multiplicative constant involving the size of the space Θ and the magnitude of the (sub)gradients of the losses ℓ_t —are typical and, as we shall see, minimax optimal.

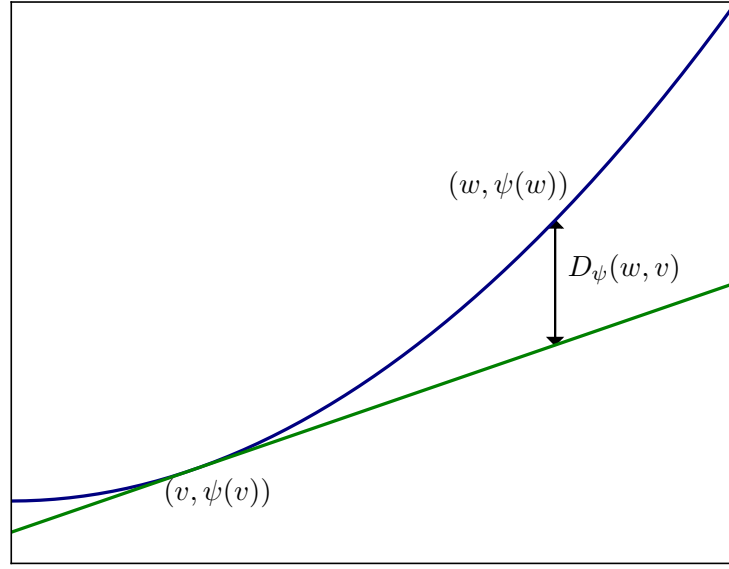


Figure 17.2: Illustration of Bregman divergence.

17.2.2 Mirror descent-type methods

In a variety of scenarios, it is advantageous to measure distances in a way more amenable to the problem structure, for example, if Θ is a probability simplex or we have prior information about the loss functions ℓ_t that nature may choose. With this in mind, we present a slightly more general algorithm than the projected gradient method 17.1, which requires us to give a few more definitions.

Given a convex differentiable function $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$, we will replace the quadratic in the update (17.2.2) with a *Bregman divergence* associated with ψ , where we recall

$$D_\psi(w, v) = \psi(w) - \psi(v) - \langle \nabla \psi(v), w - v \rangle. \quad (17.2.4)$$

The Bregman divergence is always non-negative, as $D_\psi(w, v)$ is the gap between the true function value $\psi(w)$ and its linear approximation at the point v (see Figure 17.2).

Example 17.2.3 (Euclidean distance as Bregman divergence): Take $\psi(w) = \frac{1}{2} \|w\|_2^2$ to obtain $D(w, v) = \frac{1}{2} \|w - v\|_2^2$. More generally, if for a matrix A we define $\|w\|_A^2 = w^\top A w$, then taking $\psi(w) = \frac{1}{2} w^\top A w$, we have

$$D_\psi(w, v) = \frac{1}{2} (w - v)^\top A (w - v) = \frac{1}{2} \|w - v\|_A^2.$$

So Bregman divergences generalize (squared) Euclidean distance. \diamond

Example 17.2.4 (KL divergence as a Bregman divergence): Take $\psi(w) = \sum_{j=1}^d w_j \log w_j$. Then ψ is convex over the positive orthant \mathbb{R}_+^d (the second derivative of $w \log w$ is $1/w$), and for $w, v \in \Delta_d = \{u \in \mathbb{R}_+^d : \langle \mathbf{1}, u \rangle = 1\}$, we have

$$D_\psi(w, v) = \sum_j w_j \log w_j - \sum_j v_j \log v_j - \sum_j (1 + \log v_j)(w_j - v_j) = \sum_j w_j \log \frac{w_j}{v_j} = D_{\text{kl}}(w \| v),$$

where in the final equality we treat w and v as probability distributions on $\{1, \dots, d\}$. \diamond

With these examples in mind, we now present the mirror descent algorithm, which is the natural generalization of online gradient descent. We call ψ the *distance generating function* as it yields the Bregman divergence in the method, functioning analogously to the squared ℓ_2 -distance in the earlier Algorithm 17.1.

Input: distance-generating function ψ , parameter space Θ , and non-increasing stepsize sequence η_1, η_2, \dots

Repeat: for each iteration t ,

- i. predict $\theta_t \in \Theta$, receive function ℓ_t and suffer loss $\ell_t(\theta_t)$.
- ii. construct model $\hat{\ell}_t$ of ℓ_t at θ_t , and perform non-Euclidean update

$$\theta_{t+1} = \operatorname{argmin}_{\theta \in \Theta} \left\{ \hat{\ell}_t(\theta; \theta_t) + \frac{1}{\eta_t} D_\psi(\theta, \theta_t) \right\}. \quad (17.2.5)$$

Figure 17.3: The online model-based mirror descent algorithm

The mirror descent update (17.2.5) frequently admits easy-to-compute solutions. Assume for simplicity that we use the linear model $\hat{\ell}_t(\theta) = \ell_t(\theta_t) + \langle g_t, \theta - \theta_t \rangle$, where $g_t \in \partial \ell_t(\theta_t)$. By taking $\Theta = \mathbb{R}^d$ and $\psi(\theta) = \frac{1}{2} \|\theta\|_2^2$, we note that the mirror descent procedure simply corresponds to the gradient update $\theta_{t+1} = \theta_t - \eta_t g_t$, and it evidently generalizes the update (17.2.2). We can also recover the *exponentiated gradient* or *entropic mirror descent* algorithm.

Example 17.2.5 (Exponentiated gradients): Let $\Theta = \Delta_d = \{v \in \mathbb{R}_+^d : \langle \mathbf{1}, v \rangle = 1\}$, the probability simplex in \mathbb{R}^d . Then a natural choice for ψ is the negative entropy, $\psi(w) = \sum_j w_j \log w_j$, which (as in Example 17.2.4) gives $D_\psi(w, v) = \sum_j w_j \log \frac{w_j}{v_j}$.

Consider the update step (17.2.5). Fixing $v = \theta_t$ for notational simplicity, we must solve

$$\underset{\theta}{\text{minimize}} \quad \langle g, \theta \rangle + \frac{1}{\eta} \sum_j \theta_j \log \frac{\theta_j}{v_j} \quad \text{subject to } \theta \in \Delta_d.$$

Writing the Lagrangian for this problem after introducing multipliers $\tau \in \mathbb{R}$ for the constraint that $\langle \mathbf{1}, \theta \rangle = 1$ and $\lambda \in \mathbb{R}_+^d$ for $\theta \succeq 0$, we have

$$\mathcal{L}(\theta, \lambda, \tau) = \langle g, \theta \rangle + \frac{1}{\eta} \sum_{j=1}^d \theta_j \log \frac{\theta_j}{v_j} - \langle \lambda, \theta \rangle + \tau (\langle \mathbf{1}, \theta \rangle - 1),$$

which we minimize by taking

$$\theta_j = v_j \exp(-\eta g_j + \lambda_j \eta - \tau \eta - 1).$$

As $\theta_j > 0$ certainly, the constraint $\theta \succeq 0$ is inactive and $\lambda_j = 0$. Thus, choosing τ to normalize the vector θ , we obtain the *exponentiated gradient update*

$$\theta_{t+1,i} = \frac{\theta_{t,i} e^{-\eta_t g_{t,i}}}{\sum_j \theta_{t,j} e^{-\eta_t g_{t,j}}} \quad \text{for } i = 1, \dots, d,$$

the explicit form of the update (17.2.5) when using the linear approximation for $\hat{\ell}$. \diamond

17.2.3 Convergence analysis of mirror descent

We now turn to an analysis of the mirror descent algorithm. Before presenting the analysis, we require two more definitions that allow us to relate Bregman divergences to various norms.

Definition 17.2. Let $\|\cdot\|$ be a norm. The dual norm $\|\cdot\|_*$ associated with $\|\cdot\|$ is

$$\|y\|_* := \sup_{x: \|x\| \leq 1} x^\top y.$$

For example, a straightforward calculation shows that the dual to the ℓ_∞ -norm is the ℓ_1 -norm, and the Euclidean norm $\|\cdot\|_2$ is self-dual (by the Cauchy-Schwarz inequality). Lastly, we require a definition of functions of suitable curvature for use in mirror descent methods.

Definition 17.3. A convex function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is strongly convex with respect to the norm $\|\cdot\|$ over the set Θ if for all $w, v \in \Theta$ and $g \in \partial f(w)$ we have

$$f(v) \geq f(w) + \langle g, v - w \rangle + \frac{1}{2} \|w - v\|^2.$$

The function f is strongly convex if it grows at least quadratically fast at every point in its domain. Strongly convex functions play an important role in the stability properties of minimizers and enjoy a number of equivalent characterizations; see Proposition C.1.5 in Appendix C.1 for more.

The definition (17.2.4) of the divergence makes apparent that ψ is strongly convex if and only if

$$D_\psi(w, v) \geq \frac{1}{2} \|w - v\|^2.$$

As three examples, we consider Euclidean distance, entropy, and p -norms for $1 < p \leq 2$.

Example 17.2.6: For the Euclidean distance, which uses $\psi(w) = \frac{1}{2} \|w\|_2^2$, we have $\nabla \psi(w) = w$, and

$$\frac{1}{2} \|v\|_2^2 = \frac{1}{2} \|w + v - w\|_2^2 = \frac{1}{2} \|w\|_2^2 + \langle w, v - w \rangle + \frac{1}{2} \|w - v\|_2^2$$

by a calculation, so that ψ is strongly convex with respect to the ℓ_2 -norm. \diamond

Example 17.2.7: Let $\psi(w) = \sum_j w_j \log w_j$ be the negative entropy. Then ψ is strongly convex with respect to the ℓ_1 -norm, that is,

$$D_\psi(w, v) = D_{\text{kl}}(w \| v) \geq \frac{1}{2} \|w - v\|_1^2,$$

which follows immediately from Pinsker's inequality, Proposition 2.2.8. \diamond

It can be convenient to have divergences strongly convex with respect to the ℓ_1 -norm over all of \mathbb{R}^d rather than just the positive orthant. Squared ℓ_p -norms provide this guarantee:

Example 17.2.8: Let $\psi(w) = \frac{1}{2(p-1)} \|w\|_p^2$, where $1 < p \leq 2$. Then ψ is strongly convex with respect to the ℓ_p -norm $\|\cdot\|_p$. (Exercise 17.5 asks you to prove this fact.)

In dimension $d \geq 3$, consider the choice $p = 1 + \frac{1}{\log d}$. Then because $\|w\|_1 \leq d^{\frac{p-1}{p}} \|w\|_p$ for all $p \in [1, \infty]$, then for this choice of p we have $\|w\|_p \geq d^{\frac{1-p}{p}} \|w\|_1 = d^{-\frac{1}{1+\log d}} \|w\|_1 \geq e^{-1} \|w\|_1$, so

$$D_\psi(w, v) \geq \frac{1}{2} \|w - v\|_p^2 \geq \frac{1}{2e^2} \|w - v\|_1^2.$$

To within a numerical constant, $\psi(w) = \frac{1}{2(p-1)} \|w\|_p^2$ is strongly convex with respect to the ℓ_1 -norm. \diamond

With these examples in place, we present the main convergence guarantee for mirror descent.

Theorem 17.2.9 (Regret of mirror descent). *Let ℓ_t be an arbitrary sequence of convex functions, and let θ_t be generated according to the mirror descent algorithm 17.3. Assume that the proximal function ψ is strongly convex with respect to the norm $\|\cdot\|$, which has dual norm $\|\cdot\|_*$. Then for a sequence of subgradients $g_t \in \partial\ell_t(\theta_t)$,*

(a) *If $\eta_t = \eta$ for all t , then for any $\theta^* \in \Theta$,*

$$\sum_{t=1}^n [\ell_t(\theta_t) - \ell_t(\theta^*)] \leq \frac{1}{\eta} D_\psi(\theta^*, \theta_1) + \frac{\eta}{2} \sum_{t=1}^n \|g_t\|_*^2.$$

(b) *If Θ is compact and $D_\psi(\theta^*, \theta) \leq R^2$ for any $\theta \in \Theta$, then*

$$\sum_{t=1}^n [\ell_t(\theta_t) - \ell_t(\theta^*)] \leq \frac{1}{\eta_n} R^2 + \sum_{t=1}^n \frac{\eta_t}{2} \|g_t\|_*^2.$$

We defer the proof temporarily to Section 17.2.5.

17.2.4 Instantiations of the regret guarantee

Before proving Theorem 17.2.9, we provide instantiations to exhibit its guarantees. First, we note that for Lipschitz continuity properties of the losses are equivalent to bounds on the norms of the subgradients g_t : indeed, the following two conditions are equivalent: (a) ℓ_t is G -Lipschitz with respect to the norm $\|\cdot\|$, meaning that

$$|\ell_t(\theta) - \ell_t(\theta')| \leq G \|\theta - \theta'\|$$

for all θ, θ' , and (b)

$$\|g\|_* \leq G \text{ for all } g \in \partial\ell_t(\theta) \text{ and } \theta.$$

(See Exercise 17.6.) When designing a particular instantiation of the online mirror descent algorithm 17.3, we therefore carefully consider the strong convexity properties of ψ that we employ, as these allow control over the gradient magnitudes in Theorem 17.2.9 via the dual norms $\|g_t\|_*$.

Example 17.2.10 (Lipschitz constants and linear prediction): Consider margin-based binary classification, where for data $(x_t, y_t) \in \mathbb{R}^d \times \{\pm 1\}$ we have loss $\ell_t(\theta) = \phi(y_t \langle \theta, x_t \rangle)$ for some convex ϕ with $\phi'(0) < 0$ and $|\phi(s)| \leq 1$ for any $\phi'(s) \in \partial\phi(s)$, $s \in \mathbb{R}$. Particular choices include

- i. the logistic loss, with $\phi(t) = \log(1 + e^{-t})$, and
- ii. the hinge loss, with $\phi(t) = [1 - t]_+$.

Then $\partial\ell_t(\theta) = y_t x_t \partial\phi(y_t \langle \theta, x_t \rangle)$ and the losses ℓ_t are G -Lipschitz with respect to the norm $\|\cdot\|$ if and only if $\|x\|_* \leq G$ for all $x \in \mathcal{X}$. \diamond

In the Euclidean case, when $\psi(\theta) = \frac{1}{2} \|\theta\|_2^2$, we assume that the loss functions ℓ_t are all G -Lipschitz with respect to the ℓ_2 -norm. In Example 17.2.10, this corresponds to the data x belonging to an ℓ_2 -ball of radius G . In this case, the two regret bounds above become

$$\frac{1}{2\eta} \|\theta^* - \theta_1\|_2^2 + \frac{\eta}{2} n G^2 \quad \text{and} \quad \frac{1}{2\eta_n} R^2 + \sum_{t=1}^n \frac{\eta_t}{2} G^2,$$

respectively, where in the second case we assumed that $\|\theta^* - \theta_t\|_2 \leq R$ for all t . In the former case, we take $\eta = \frac{R}{G\sqrt{n}}$ and recover Corollary 17.2.2; in the latter we take $\eta_t = \frac{R}{G\sqrt{t}}$, which does not require knowledge of n ahead of time. We obtain the following corollary.

Corollary 17.2.11. *Assume that $\Theta \subset \{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq R\}$ and that the loss functions ℓ_t are G -Lipschitz with respect to the Euclidean norm. Take $\eta_t = \frac{R}{G\sqrt{t}}$. Then for all $\theta^* \in \Theta$,*

$$\sum_{t=1}^n [\ell_t(\theta_t) - \ell_t(\theta^*)] \leq 3RG\sqrt{n}.$$

Proof For any $\theta, \theta^* \in \Theta$, we have $\|\theta - \theta^*\|_2 \leq 2R$, so that $D_\psi(\theta^*, \theta) \leq 4R^2$. Using that

$$\sum_{t=1}^n t^{-\frac{1}{2}} \leq \int_0^n t^{-\frac{1}{2}} dt = 2\sqrt{n}$$

gives the result. \square

Other geometries suggest alternative choices of the distance-generating function ψ . For example, in high-dimensional settings, when the underlying domain Θ is the probability simplex, we can achieve bounds that depend only on the ℓ_∞ norm of the gradients. (Recall Example 17.1.3 for motivation.) In this case, we have the following corollary to Theorem 17.2.9.

Corollary 17.2.12. *Assume that $\Theta = \Delta_d = \{\theta \in \mathbb{R}_+^d : \langle \mathbf{1}, \theta \rangle = 1\}$ and take the proximal function $\psi(\theta) = \sum_j \theta_j \log \theta_j$ to be the negative entropy in the mirror descent procedure 17.3. Then with the fixed stepsize η and initial point as the uniform distribution $\theta_1 = \mathbf{1}/d$, we have for any sequence of convex losses ℓ_t*

$$\sum_{t=1}^n [\ell_t(\theta_t) - \ell_t(\theta^*)] \leq \frac{\log d}{\eta} + \frac{\eta}{2} \sum_{t=1}^n \|g_t\|_\infty^2.$$

Proof Using Pinsker's inequality in the form of Example 17.2.7, we have that ψ is strongly convex with respect to $\|\cdot\|_1$. Consequently, taking the dual norm to be the ℓ_∞ -norm, part (a) of Theorem 17.2.9 shows that

$$\sum_{t=1}^n [\ell_t(w_t) - \ell_t(w^*)] \leq \frac{1}{\eta} \sum_{j=1}^d w_j^* \log \frac{w_j^*}{w_{1,j}} + \frac{\eta}{2} \sum_{t=1}^n \|g_t\|_\infty^2.$$

Noting that with $w_1 = \mathbf{1}/d$, we have $D_\psi(w^*, w_1) \leq \log d$ for any $w^* \in \Theta$ gives the result. \square

Corollary 17.2.12 yields somewhat sharper results than Corollary 17.2.11, though in the restricted setting that Θ is the probability simplex in \mathbb{R}^d . Indeed, let us assume that the subgradients $g_t \in [-1, 1]^d$, the hypercube in \mathbb{R}^d . In this case, the tightest possible bound on their ℓ_2 -norm is $\|g_t\|_2 \leq \sqrt{d}$, while $\|g_t\|_\infty \leq 1$ always. Similarly, if $\Theta = \Delta_d$, then we may only guarantee that $\|\theta^* - \theta_1\|_2 \leq \sqrt{2}$. Thus, Euclidean methods (Corollary 17.2.11) can only guarantee regret

$$\frac{1}{2\eta} \|\theta^* - \theta_1\|_2^2 + \frac{\eta}{2} nd \leq \sqrt{2nd} \text{ with the choice } \eta = \frac{1}{\sqrt{2nd}},$$

while the entropic mirror descent procedure (Alg. 17.3 with $\psi(\theta) = \sum_j \theta_j \log \theta_j$) guarantees

$$\frac{\log d}{\eta} + \frac{\eta}{2}n \leq \sqrt{2n \log d} \quad \text{with the choice } \eta = \frac{\sqrt{2 \log d}}{2\sqrt{n}}. \quad (17.2.6)$$

The latter guarantee is *exponentially* better in the dimension.

As a final example, we revisit the p -norm algorithms in high dimensions (at least, when $d \geq 3$).

Example 17.2.13 (p -norm algorithms): Letting $1 < p \leq 2$, so that $\psi(\theta) = \frac{1}{2(p-1)} \|\theta\|_p^2$ is strongly convex with respect to $\|\cdot\|_p$, as in Example 17.2.8 over \mathbb{R}^d . Exercise 17.8 explores the computation of the update (17.2.5), so we focus exclusively on the regret guarantee. Assume w.l.o.g. that $\mathbf{0} \in \Theta$, and let $\theta_1 = \mathbf{0}$ (as otherwise, we take $\psi(\theta) = \frac{1}{2(p-1)} \|\theta - \theta_1\|_p^2$, and the same results hold). Consider the fixed stepsize result in Theorem 17.2.9. Then for $q = \frac{p}{p-1}$ we obtain

$$\sum_{t=1}^n \ell_t(\theta_t) - \ell_t(\theta^*) \leq \frac{\|\theta^*\|_p^2}{2(p-1)\eta} + \frac{\eta}{2} \sum_{t=1}^n \|g_t\|_q^2.$$

Choose $p = 1 + \frac{1}{\log d} < 2$, so that its conjugate $q = \frac{p}{p-1} = 1 + \log d$ and satisfies $\|g\|_q \leq d^{1/q} \|g\|_\infty \leq e \|g\|_\infty$. Then because $\|\theta\|_p \leq \|\theta\|_1$, we have

$$\sum_{t=1}^n [\ell_t(\theta_t) - \ell_t(\theta^*)] \leq \frac{\log d}{2\eta} \|\theta^*\|_1^2 + \frac{e\eta}{2} \sum_{t=1}^n \|g_t\|_\infty^2.$$

Assume each ℓ_t is G_∞ -Lipschitz with respect to the ℓ_1 -norm (which is easier to satisfy than Lipschitzness with respect to any other ℓ_p -norm, as it is equivalent to boundedness of the gradients in ℓ_∞), and that the domain Θ is contained in an ℓ_1 -ball of ℓ_1 -radius R_1 . Then with the choice $\eta = R_1 \sqrt{\log d} / G_\infty \sqrt{n}$ we obtain

$$\sum_{t=1}^n [\ell_t(\theta_t) - \ell_t(\theta^*)] \leq O(1) G_\infty R_1 \sqrt{n \log d}.$$

As in Corollary 17.2.12, we see the improvement in dimension dependence over Euclidean-type methods. \diamond

17.2.5 Proof of Theorem 17.2.9

The proof of the theorem proceeds in three lemmas, which distill various optimality conditions for convex optimization problems, and then an inductive application.

The first (see Proposition C.1.4 in Appendix C.1) explicitly characterizes optimality for a convex optimization problem.

Lemma 17.2.14. *Let $h : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function and Θ be a convex set. Then θ^* minimizes h over Θ if and only if there exists $g \in \partial h(\theta^*)$ such that*

$$\langle g, \theta - \theta^* \rangle \geq 0 \quad \text{for all } \theta \in \Theta.$$

The next result gives the key progress guarantee for any model-based minimization strategy.

Lemma 17.2.15. *Let $\hat{\ell}$ be any valid model (Definition 17.1) of $\ell : \Theta \rightarrow \mathbb{R}$ at the point θ_0 , and for $\eta > 0$ define*

$$\theta_\eta := \operatorname{argmin}_{\theta \in \Theta} \left\{ \hat{\ell}(\theta) + \frac{1}{\eta} D_\psi(\theta, \theta_0) \right\}.$$

Then for some $g_0 \in \partial \ell(\theta_0)$ and any $\theta \in \Theta$,

$$\ell(\theta_0) - \ell(\theta) \leq \frac{1}{\eta} [D_\psi(\theta, \theta_0) - D_\psi(\theta, \theta_\eta) - D_\psi(\theta_\eta, \theta_0)] + \langle g_0, \theta_0 - \theta \rangle.$$

Proof We consider the optimality conditions for the minimization. By the optimality conditions in Lemma 17.2.14, there exists $g_\eta \in \partial \hat{\ell}(\theta_\eta)$ such that

$$\langle g_\eta + \eta^{-1}(\nabla \psi(\theta_\eta) - \nabla \psi(\theta_0)), \theta - \theta_\eta \rangle \geq 0 \quad \text{for all } \theta \in \Theta.$$

Rearranging, we have

$$\frac{1}{\eta} \langle \nabla \psi(\theta_\eta) - \nabla \psi(\theta_0), \theta - \theta_\eta \rangle \geq \langle g_\eta, \theta_\eta - \theta \rangle.$$

The magical step is the “three term identity,” valid for any u, v, w , that

$$D_\psi(u, v) - D_\psi(u, w) - D_\psi(w, v) = \langle \nabla \psi(w) - \nabla \psi(v), u - w \rangle,$$

and substituting $u = \theta$, $v = \theta_0$, and $w = \theta_\eta$, we find that

$$\langle g_\eta, \theta_\eta - \theta \rangle \leq \frac{1}{\eta} [D_\psi(\theta, \theta_0) - D_\psi(\theta, \theta_\eta) - D_\psi(\theta_\eta, \theta_0)]. \quad (17.2.7)$$

Finally, we use the modeling conditions in Definition 17.1 to lower bound the left hand side. To that end, note that if $g_0 \in \partial \ell(\theta_0)$, then by convexity we obtain

$$\begin{aligned} \langle g_\eta, \theta_\eta - \theta \rangle &\stackrel{(i)}{\geq} \hat{\ell}(\theta_\eta) - \hat{\ell}(\theta) \stackrel{(ii)}{\geq} \hat{\ell}(\theta_\eta) - \ell(\theta) \stackrel{(iii)}{\geq} \hat{\ell}(\theta_0) + \langle g_0, \theta_\eta - \theta_0 \rangle - \ell(\theta) \\ &= \ell(\theta_0) + \langle g_0, \theta_\eta - \theta_0 \rangle - \ell(\theta) \end{aligned}$$

where inequalities (i) and (iii) are first-order convexity, and inequality (ii) follows from the lower bounding condition in Definition 17.1. Substituting and rearranging in inequality (17.2.7) yields the claimed inequality in the lemma, and the last bit is to use the inclusion (17.2.1), that if $g_0 \in \partial \ell(\theta_0)$ then $g_0 \in \partial \ell(\theta_0)$. \square

Summing the inequality in Lemma 17.2.15 implies the following regret bound.

Lemma 17.2.16. *Let $\ell_t : \Theta \rightarrow \mathbb{R}$ be any sequence of convex loss functions and η_t be a non-increasing sequence, where $\eta_0 = \infty$. Then with the mirror descent strategy (17.2.5), there are subgradient $g_t \in \partial \ell_t(\theta_t)$ such that for any $\theta^* \in \Theta$ we have*

$$\sum_{t=1}^n \ell_t(\theta_t) - \ell_t(\theta^*) \leq \sum_{t=1}^n \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) D_\psi(\theta^*, \theta_t) + \sum_{t=1}^n \left[-\frac{1}{\eta_t} D_\psi(\theta_{t+1}, \theta_t) + \langle g_t, \theta_t - \theta_{t+1} \rangle \right].$$

Proof Summing the inequality in Lemma 17.2.15 gives

$$\begin{aligned} \sum_{t=1}^n \ell_t(\theta_t) - \ell_t(\theta^*) &\leq \sum_{t=1}^n \frac{1}{\eta_t} [D_\psi(\theta^*, \theta_t) - D_\psi(\theta^*, \theta_{t+1}) - D_\psi(\theta_{t+1}, \theta_t)] + \sum_{t=1}^n \langle g_t, \theta_t - \theta_{t+1} \rangle \\ &= \sum_{t=2}^n \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) D_\psi(\theta^*, \theta_t) + \frac{1}{\eta_1} D_\psi(\theta^*, \theta_1) - \frac{1}{\eta_n} D_\psi(\theta^*, \theta_{n+1}) \\ &\quad + \sum_{t=1}^n \left[-\frac{1}{\eta_t} D_\psi(\theta_{t+1}, \theta_t) + \langle g_t, \theta_t - \theta_{t+1} \rangle \right] \end{aligned}$$

as desired. \square

It remains to use the negative terms $-D_\psi(\theta_t, \theta_{t+1})$ to cancel the gradient terms $\langle g_t, \theta_t - \theta_{t+1} \rangle$. To that end, we recall Definition 17.2 of the dual norm $\|\cdot\|_*$ and the strong convexity assumption on ψ . Using the Fenchel-Young inequality that $ab \leq \frac{1}{2\eta} a^2 + \frac{\eta}{2} b^2$, valid for any $\eta > 0$, we have

$$\langle g_t, \theta_t - \theta_{t+1} \rangle \leq \|g_t\|_* \|\theta_t - \theta_{t+1}\| \leq \frac{\eta_t}{2} \|g_t\|_*^2 + \frac{1}{2\eta_t} \|\theta_t - \theta_{t+1}\|^2.$$

Now, we use the strong convexity condition, which gives

$$-\frac{1}{\eta_t} D_\psi(\theta_{t+1}, \theta_t) \leq -\frac{1}{2\eta_t} \|\theta_t - \theta_{t+1}\|^2.$$

Combining the preceding two displays in Lemma 17.2.16 gives the result of Theorem 17.2.9.

17.3 Optimality guarantees and fundamental limits

Developing minimax and other types of lower bounds for stochastic optimization and general statistical learning problems requires some additional technology beyond the minimax bounds we have already developed, as we no longer presume there is some parameter θ to actually estimate. Instead, we only wish to achieve small expected loss $L(\theta) = \mathbb{E}[\ell(\theta; Z)]$ in the stochastic optimization case (17.0.2), or small regret (17.0.1). Thus, we redefine the typical minimax risk to consider only the gap in the population losses, so that for any loss function $\ell : \Theta \times \mathcal{Z} \rightarrow \mathbb{R}$, we use the shorthand $L_P(\theta) = \mathbb{E}_P[\ell(\theta, Z)]$ and define the minimax (optimization) risk

$$\mathfrak{M}_n(\mathcal{P}, \Theta, \ell) := \inf_{\hat{\theta} : \mathcal{Z}^n \rightarrow \Theta} \sup_{P \in \mathcal{P}} \left\{ \mathbb{E}_{P^n} \left[L_P(\hat{\theta}(Z_1^n)) \right] - \inf_{\theta \in \Theta} L_P(\theta) \right\}. \quad (17.3.1)$$

Considering the convergence guarantees in the preceding section, it is also interesting to provide minimax lower bounds that capture the geometric aspects of the problems we consider—their Lipschitz continuity properties and relationship with the underlying domain $\Theta \subset \mathbb{R}^d$. Because it is tedious to carefully address the interaction between probability distributions and losses ℓ , we modify the minimax definition (17.3.1) slightly. Now, let \mathcal{L} consist of functions $\ell : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$. Let $\mathcal{P}(\mathcal{L})$ consist of all probability distributions on $\ell \in \mathcal{L}$, so that we identify the sample space \mathcal{Z} with losses themselves, and make the natural definition $L_P(\theta) = \mathbb{E}_P[\ell(\theta)] = \int \ell(\theta) dP(\ell)$. This

makes it easy, for example, to consider the class of convex functions that are Lipschitz continuous with respect to the p -norm

$$\mathcal{L}_p := \left\{ \text{convex } \ell : \mathbb{R}^d \rightarrow \mathbb{R}_+ \text{ s.t. } |\ell(\theta) - \ell(\theta')| \leq \|\theta - \theta'\|_p \text{ for all } \theta, \theta' \in \mathbb{R}^d \right\}.$$

Then we let

$$\mathfrak{M}_n(\Theta, \mathcal{L}) := \inf_{\hat{\theta}_n} \sup_{P \in \mathcal{P}(\mathcal{L})} \left\{ \mathbb{E}_{P^n} [L_P(\hat{\theta}_n)] - \inf_{\theta \in \Theta} L_P(\theta) \right\}, \quad (17.3.2)$$

where the infimum is over estimators $\hat{\theta}_n$ observing n observations i.i.d. from P .

17.3.1 From optimization to testing

Because we no longer measure parameter error, we introduce a more general lower bounding technique, which employs the same general idea of our reductions from estimation to testing in Chapter 9. Now, however, we wish to reduce optimization to testing, constructing problems where optimizing well implies that we can solve certain hypothesis tests. Thus, we define new versions of separations (analogues of packings), and show how we can use them to show that optimizing well implies accurately testing between different probability distributions.

We begin by defining separations in terms of optimality gaps. For a parameter space Θ and collection \mathcal{P} of distributions, a function

$$r : \Theta \times \mathcal{P} \rightarrow \mathbb{R}_+$$

is a valid *risk gap* if

$$\inf_{\theta \in \Theta} r(\theta, P) = 0 \text{ for each } P \in \mathcal{P}. \quad (17.3.3)$$

Example 17.3.1 (A risk gap in stochastic optimization): Stochastic optimization problems immediately suggest a particular choice for the risk gap: given any loss $\ell : \Theta \times \mathcal{Z} \rightarrow \mathbb{R}$ and associated population losses $L_P(\theta) := \mathbb{E}_P[\ell(\theta, Z)]$, the functional

$$r(\theta, P) := L_P(\theta) - \inf_{\theta \in \Theta} L_P(\theta) \quad (17.3.4)$$

immediately satisfies the condition (17.3.3). \diamond

JCD Comment: Make a figure of separation.

We now mimic the reduction from estimation to testing we developed in Chapter 9.2.1. With this in mind, for distributions P_0, P_1 , we define the *separation* between them (for the risk gap r) by

$$\text{sep}_r(P_0, P_1; \Theta) := \sup \left\{ \delta \geq 0 : \begin{array}{l} r(\theta, P_0) \leq \delta \text{ implies } r(\theta, P_1) \geq \delta \\ r(\theta, P_1) \leq \delta \text{ implies } r(\theta, P_0) \geq \delta \end{array} \text{ for any } \theta \in \Theta \right\}. \quad (17.3.5)$$

That is, having small loss on P_0 implies large loss on P_1 and vice versa. The next examples show simple ways to construct separated objectives in optimization using absolute and squared losses. Note that they have different scaling in the natural underlying “separation,” which suggests—correctly—that optimizing one loss may be easier than the other.

Example 17.3.2 (A separation for one-dimensional objectives): Consider data $x \in \{-1, 1\}$, and losses $\ell(\theta, x) = |\theta - x|$. Consider the distributions P_{-1} and P_1 on X with $P_v(X = v) = \frac{1+\delta}{2}$ and $P_v(X = -v) = \frac{1-\delta}{2}$ for $v \in \{-1, 1\}$, where $\delta \in [0, 1]$. Then by inspection,

$$L_v(\theta) := \mathbb{E}_{P_v}[\ell(\theta, X)] = \frac{1+\delta}{2}|\theta - v| + \frac{1-\delta}{2}|\theta + v|$$

satisfies $\operatorname{argmin}_{\theta} L_v(\theta) = v$, and $L_v^* := \inf_{\theta} L_v(\theta) = (1 - \delta)$. For the risk gap (17.3.4), we use that if $\theta v \leq 0$ then $L_v(\theta) - L_v^* \geq 1 - (1 - \delta) = \delta$, so that

$$\operatorname{sep}_r(P_{-1}, P_1; \mathbb{R}) = \operatorname{sep}_r(P_{-1}, P_1; [-1, 1]) = \delta.$$

We have a separation for any domain $\Theta \supset [-1, 1]$. \diamond

Example 17.3.3 (Separation of quadratic losses): As in the preceding example, consider data $x \in \{-1, 1\}$, but take squared losses $\ell(\theta, x) = \frac{1}{2}(\theta - x)^2$. Then for Bernoulli distributions $P_v(X = v) = \frac{1+\delta}{2}$ for $v \in \{-1, 1\}$ and $\delta \in [0, 1]$, we have $L_v(\theta) := \mathbb{E}_{P_v}[\ell(\theta, X)] = \frac{1+\delta}{2}(\theta - v)^2 + \frac{1-\delta}{2}(\theta + v)^2$. For these expected losses, we have

$$\operatorname{argmin}_{\theta} L_v(\theta) = \delta v \quad \text{and} \quad \inf_{\theta} L_v(\theta) = \frac{1+\delta}{2}(1 - \delta)^2 + \frac{1-\delta}{2}(1 + \delta)^2 = 1 - \delta^2.$$

In this case, for $\theta v \leq 0$ we have $L_v(\theta) \geq 1$ and so

$$\operatorname{sep}_r(P_{-1}, P_1; \mathbb{R}) = \operatorname{sep}_r(P_{-1}, P_1; [-1, 1]) = \delta^2,$$

a quadratically smaller separation than for the absolute loss in Example 17.3.2. \diamond

To complete our reduction of optimization to testing, we require separation between collections of distributions. Thus, for a collection of distributions $\{P_v\}_{v \in \mathcal{V}}$ indexed by \mathcal{V} , we say $\{P_v\}$ is δ -separated if

$$\operatorname{sep}_r(P_v, P_{v'}; \Theta) \geq \delta \quad \text{for } v \neq v' \in \mathcal{V}.$$

Then by an argument similar to Proposition 9.2.1, we have the following reduction from *optimization* to testing.

Proposition 17.3.4. *Let $\{P_v\}_{v \in \mathcal{V}} \subset \mathcal{P}$ be δ -separated for the risk gap $r : \Theta \times \mathcal{P} \rightarrow \mathbb{R}_+$. Then for any estimator $\hat{\theta} : \mathcal{Z} \rightarrow \Theta$,*

$$\frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \mathbb{E}_{P_v}[r(\hat{\theta}, P_v)] \geq \delta \inf_{\hat{v}} \mathbb{P}(\hat{v} \neq V), \quad (17.3.6)$$

where \mathbb{P} is the joint distribution over the random index V chosen uniformly and, conditional on $V = v$, drawing $Z \sim P_v$. If additionally $r(\theta, P) \geq 0$ for all θ , inequality (17.3.6) holds for any estimator rather than estimators taking values only in Θ .

Proof Suppose that $\{P_v\}$ is δ -separated. Then

$$\frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \mathbb{E}_{P_v}[r(\hat{\theta}, P_v)] \geq \delta \mathbb{P}(r(\hat{\theta}, P_V) \geq \delta).$$

Define $\Psi : \Theta \rightarrow \mathcal{V}$ by

$$\Psi(\theta) = \begin{cases} v & \text{if } r(\theta, P_v) < \delta \\ \text{arbitrary} & \text{otherwise.} \end{cases}$$

Then if $r(\hat{\theta}, P_v) < \delta$, we must have $\Psi(\theta) = v$, and so $\mathbb{P}(r(\hat{\theta}, P_V) < \delta) \leq \mathbb{P}(\Psi(\hat{\theta}) = V)$, that is,

$$\mathbb{P}(r(\hat{\theta}, P_V) \geq \delta) = 1 - \mathbb{P}(r(\hat{\theta}, P_V) < \delta) \geq 1 - \mathbb{P}(\Psi(\hat{\theta}) = V) = \mathbb{P}(\Psi(\hat{\theta}) \neq V).$$

Take an infimum over tests to get the result. \square

As a corollary, we obtain a lower bound on the minimax risk for stochastic optimization. As we have already seen in Example 17.3.1, $r(\theta, P) = L_P(\theta) - \inf_{\theta^* \in \Theta} L_P(\theta^*)$ is a valid risk gap, making the following result immediate.

Corollary 17.3.5. *Assume that \mathcal{P} has a δ -separated subset $\{P_v\}_{v \in \mathcal{V}}$ for the population loss gap (17.3.4). Then the minimax risk (17.3.1) satisfies*

$$\mathfrak{M}_n(\mathcal{P}, \Theta, \ell) \geq \delta \inf_{\hat{v}} \mathbb{P}(\hat{v}(Z_1^n) \neq V),$$

where \mathbb{P} denotes the joint distribution over V following an arbitrary distribution and, conditional on $V = v$, drawing $Z_1^n \stackrel{\text{iid}}{\sim} P_v$.

17.3.2 Constructing hard classes of optimization problems

Given Corollary 17.3.5, our program is clear: for a loss ℓ , we construct a well-separated collection of distributions $\{P_v\}$ —where the separation scales with some parameter δ —then show that the information the observations $Z_1^n \stackrel{\text{iid}}{\sim} P_v$ contain about the particular element $v \in \mathcal{V}$ is limited. Then we apply one of our standard tools—the Assouad, Fano or Le Cam method—to lower bound the testing error, and choose $\delta > 0$ as large as possible while keeping the probability of error $\mathbb{P}(\hat{v} \neq V)$ a constant. In this section, we follow this program and provide constructions that allow us to provide explicit and optimal (to within numerical constants) lower bounds for the minimax risk for optimization problems over domains containing scaled ℓ_p balls, where the subgradients of the losses belong to $\ell_{p'}$ balls (which need not necessarily be dual to one another).

The separation (17.3.5) can be hard to compute explicitly, so we introduce an alternative that provides a sometimes simpler lower bound on the separation. For any two functions L_0, L_1 , let $\theta^v \in \operatorname{argmin}_{\theta \in \Theta} L_v(\theta)$, and define the *optimization distance* between L_0 and L_1 by

$$d_{\text{opt}}(L_0, L_1; \Theta) := \inf_{\theta \in \Theta} \{L_0(\theta) + L_1(\theta) - L_0(\theta^0) - L_1(\theta^1)\}. \quad (17.3.7)$$

This quantity always lower bounds the separation in the risk gap we use for stochastic optimization (recall Example 17.3.1).

Lemma 17.3.6. *Let P_0, P_1 be distributions and $L_v(\theta) = \mathbb{E}_{P_v}[\ell(\theta, Z)]$ for $v \in \{0, 1\}$. Then for the risk gap $r(\theta, P) := L_P(\theta) - \inf_{\theta \in \Theta} L_P(\theta)$,*

$$\operatorname{sep}_r(P_0, P_1; \Theta) \geq \frac{1}{2} d_{\text{opt}}(L_0, L_1; \Theta).$$

Proof Let $L_v^* = \inf_{\theta \in \Theta} L_v(\theta)$ for shorthand. Let $\delta \geq 0$ be such that $d_{\text{opt}}(L_0, L_1; \Theta) \geq \delta$, and assume that $L_0(\theta) - L_0^* \leq \delta/2$. Then

$$L_1(\theta) - L_1^* = L_0(\theta) + L_1(\theta) - L_0^* - L_1^* - (L_0(\theta) - L_0^*) \geq d_{\text{opt}}(L_0, L_1; \Theta) - \delta/2 \geq \delta/2.$$

The symmetric case with $L_1(\theta) - L_1^*$ is similar, so $\text{sep}_L(P_0, P_1, \Theta) \geq \frac{\delta}{2}$. \square

We now turn to the first optimality lower bound. We focus on a particular collection of distributions \mathcal{P} indexed by $v \in \{-1, 1\}^d$ on the sample space $\mathcal{X} = \{\pm e_j\}_{j=1}^d$ of signed standard basis vectors and a particular loss ℓ . Appropriate scaling will allow us to derive further lower bounds for more general parameter spaces. We follow our now familiar recipe: construct well-separated losses, obtain a packing, and then show that testing between the data coming from any particular loss is challenging.

Constructing a well-separated family of losses For $\theta \in \mathbb{R}^d$, define the loss

$$\ell(\theta; x) := \begin{cases} |\theta_j - 1| & \text{if } x = e_j \\ |\theta_j + 1| & \text{if } x = -e_j. \end{cases} \quad (17.3.8)$$

Let $v \in \{-1, 1\}^d$. For some $\delta > 0$ to be chosen, define the distribution P_v on $X \in \mathcal{X}$ by

$$X = \begin{cases} v_j e_j & \text{w.p. } \frac{1+\delta}{2d} \\ -v_j e_j & \text{w.p. } \frac{1-\delta}{2d}. \end{cases} \quad (17.3.9)$$

Then the expected loss has the explicit form

$$L_v(\theta) := \mathbb{E}_{P_v}[\ell(\theta, X)] = \frac{1}{d} \sum_{j=1}^d \left[\frac{1+\delta}{2} |\theta_j - v_j| + \frac{1-\delta}{2} |\theta_j + v_j| \right],$$

and by inspection

$$\theta_v := \underset{\theta}{\operatorname{argmin}} L_v(\theta) = v \quad \text{and} \quad \inf_{\theta} L_v(\theta) = 1 - \delta.$$

As the distance between v and v' grows, so to does the optimization distance between L_v and $L_{v'}$ for the loss (17.3.8) (as does the separation):

Lemma 17.3.7. *Let $\Theta \supset [-1, 1]^d$. For any $v, v' \in \{\pm 1\}^d$,*

$$d_{\text{opt}}(L_v, L_{v'}; \Theta) = \frac{\delta}{d} \|v - v'\|_1.$$

Proof Let $L_v^* = \inf_{\theta} L_v(\theta) = 1 - \delta$. Then writing out the quantity inside the infimum in the definition (17.3.7), we have

$$\begin{aligned} & L_v(\theta) + L_{v'}(\theta) - L_v^* - L_{v'}^* \\ &= \frac{1}{d} \sum_{j: v_j \neq v'_j} [|\theta_j - 1| + |\theta_j + 1|] + \frac{1}{d} \sum_{j: v_j = v'_j} [(1+\delta)|\theta_j - v_j| + (1-\delta)|\theta_j + v_j|] - 2(1-\delta). \end{aligned}$$

Taking infima we thus obtain

$$\inf_{\theta} \{L_v(\theta) + L_{v'}(\theta) - L_v^* - L_{v'}^*\} = \frac{1}{d} \|v - v'\|_1 + \frac{1-\delta}{d} (2d - \|v - v'\|_1) - 2(1-\delta) = \frac{\delta}{d} \|v - v'\|_1$$

as claimed. \square

Constructing the packing The actual packing construction is now straightforward: as we saw in Chapter 9.4.1, we may use the Gilbert-Varshamov bound (Lemma 9.2.3) to pack the hypercube, and then we leverage Lemma 17.3.7. Indeed, Lemma 9.2.3 implies that there exists a $d/2$ -packing of the hypercube $\{-1, 1\}^d$ in ℓ_1 distance with cardinality at least $\exp(d/8)$. We conclude the following:

Observation 17.3.8. *Let $d \in \mathbb{N}$ and $\Theta \supset [-1, 1]^d$. There exists a packing $\mathcal{V} \subset \{-1, 1\}^d$ of cardinality at least $\exp(d/8)$ such that for the loss (17.3.8) and any choice $\delta \in [0, 1]$ in the distribution (17.3.9),*

$$d_{\text{opt}}(L_v, L_{v'}; \Theta) \geq \frac{\delta}{2} \text{ for } v \neq v' \in \mathcal{V}.$$

Lower bounding the testing error Now that Observation 17.3.8 gives a separation of the expected losses (17.3.8), we can apply Fano's or Le Cam's methods to achieve a minimax lower bound on the expected optimization error. The key insights are the following KL-divergence and information bounds.

Lemma 17.3.9. *Let $X \sim P_v$ according to the coordinate sampling scheme (17.3.9). For $\delta \leq \frac{1}{2}$,*

$$D_{\text{kl}}(P_v^n \| P_{v'}^n) \leq \frac{9}{8d} \|v - v'\|_1 \delta^2.$$

For any packing $\mathcal{V} \subset \{-1, 1\}^d$, if $V \sim \text{Uniform}(\mathcal{V})$ and conditional on $V = v$ we draw $X_1^n \stackrel{\text{iid}}{\sim} P_v$,

$$I(X_1^n; V) \leq \frac{9}{4} \delta^2.$$

Proof We demonstrate the first inequality; the second follows from the trivial observation that $I(X_1^n; V) \leq \max_{v, v'} D_{\text{kl}}(P_v^n \| P_{v'}^n)$. For any $v, v' \in \{-1, 1\}^d$, the sampling scheme (17.3.9) gives

$$\begin{aligned} D_{\text{kl}}(P_v \| P_{v'}) &= \frac{1}{2d} \sum_{j: v_j \neq v'_j} \left[(1 + \delta) \log \frac{1 + \delta}{1 - \delta} + (1 - \delta) \log \frac{1 - \delta}{1 + \delta} \right] \\ &= \frac{\|v - v'\|_1}{4d} \left[(1 + \delta) \log \frac{1 + \delta}{1 - \delta} + (1 - \delta) \log \frac{1 - \delta}{1 + \delta} \right] \leq \frac{\|v - v'\|_1}{4d} (4\delta^2 + \delta^3), \end{aligned}$$

the final inequality valid for $0 \leq \delta \leq \frac{1}{2}$. Because $D_{\text{kl}}(P_v^n \| P_{v'}^n) = n D_{\text{kl}}(P_v \| P_{v'})$, we obtain $D_{\text{kl}}(P_v^n \| P_{v'}^n) \leq \frac{9}{8d} \|v - v'\|_1 \delta^2$ for $\delta \leq \frac{1}{2}$. \square

Putting everything together then yields our main lower bound.

Theorem 17.3.10. *Let $\Theta \supset [-1, 1]^d$ and \mathcal{P} contain distributions on $\{\pm e_j\}_{j=1}^d$. Then for the loss ℓ defined in (17.3.8), there exists a numerical constant $c > 0$ such that*

$$\mathfrak{M}_n(\mathcal{P}, \Theta, \ell) \geq c \min \left\{ \frac{\sqrt{d}}{\sqrt{n}}, 1 \right\}.$$

Proof Let $\mathfrak{M}_n = \mathfrak{M}_n(\mathcal{P}, \Theta, \ell)$ for simplicity, and fix $\delta \in [0, \frac{1}{4}]$ to be chosen. We consider two cases: that $d \geq 14$ and $d < 14$. In the former case, take the packing $\mathcal{V} \subset \{-1, 1\}^d$ and

distribution (17.3.9) that Observation 17.3.8 promises, which gives separation $\frac{\delta}{4}$ via Lemma 17.3.6. Then Corollary 17.3.5 and Fano's inequality imply that

$$\mathfrak{M}_n \geq \frac{\delta}{4} \left(1 - \frac{I(X_1^n; V) + \log 2}{d/8} \right).$$

By Lemma 17.3.9, we obtain

$$\mathfrak{M}_n \geq \frac{\delta}{4} \left(1 - \frac{8 \log 2}{d} - \frac{36n\delta^2}{d} \right) \geq \frac{\delta}{4} \left(\frac{3}{5} - \frac{36n\delta^2}{d} \right),$$

where we used that $d \geq 14$. Choose $\delta^2 = \min\{\frac{1}{5} \cdot \frac{d}{36n}, \frac{1}{4}\}$ to obtain the theorem when $d \geq 14$.

When $d < 14$, we simply take the packing $\mathcal{V} = \{-\mathbf{1}_d, \mathbf{1}_d\}$, which by Lemma 17.3.7 implies that $d_{\text{opt}}(L_v, L_{v'}; \Theta) = 2\delta$. Then applying Le Cam's method and Corollary 17.3.5, we have

$$\mathfrak{M}_n \geq \frac{\delta}{2} (1 - \|P_1^n - P_{-1}^n\|_{\text{TV}}) \geq \frac{\delta}{2} \left(1 - \sqrt{\frac{9n}{4}\delta^2} \right) \quad (17.3.10)$$

by Pinsker's inequality and the information bound Lemma 17.3.9 provides. Choose $\delta^2 = \min\{\frac{1}{4}, \frac{9}{16n}\}$ and treat $d < 14$ as a numerical constant. \square

17.3.3 Instantiations and optimality

By rescaling the sampling distributions and domain Θ , we can obtain minimax lower bounds for stochastic optimization with different Lipschitz constants and underlying parameter spaces. We begin by a simple rescaling.

Corollary 17.3.11. *Let the collection of losses \mathcal{L} consist of all losses G_1 -Lipschitz with respect to the ℓ_∞ -norm, meaning $g \in \partial\ell(\theta)$ implies $\|g\|_1 \leq G_1$, and assume $\Theta \supset [-R_\infty, R_\infty]$. Then*

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq cR_\infty G_1 \min \left\{ \frac{\sqrt{d}}{\sqrt{n}}, 1 \right\}.$$

Proof Scale the losses (17.3.8) and sampling distribution (17.3.9) by setting

$$X = \begin{cases} R_\infty v_j e_j & \text{w.p. } \frac{1+\delta}{2d} \\ -R_\infty v_j e_j & \text{w.p. } \frac{1-\delta}{2d} \end{cases}$$

under P_v and

$$\ell(\theta; x) = \begin{cases} G_1 |\theta_j - R_\infty| & \text{if } X_j > 0 \\ G_1 |\theta_j + R_\infty| & \text{if } X_j < 0. \end{cases}$$

The proof of Theorem 17.3.10 then proceeds *mutatis mutandis*. \square

The main foci of in our development of stochastic gradient- and mirror-descent-type methods was to develop algorithms that enjoyed convergence guarantees irrespective of the particular loss

ℓ , requiring only boundedness conditions on Θ and the gradients $g \in \partial\ell$. We can extend Theorem 17.3.10 to address this as well. For constants $C > 0$, let $C \cdot \mathcal{L} = \{C \cdot \ell \mid \ell \in \mathcal{L}\}$ denote the multiplicative scaling of a collection of losses (so, e.g., if \mathcal{L} consists of 1-Lipschitz functions, then $C \cdot \mathcal{L}$ consists of C -Lipschitz functions). The following corollary provides minimax lower bounds for such collections and domains Θ other than the ℓ_∞ -ball.

Corollary 17.3.12. *Let $2 \leq p, q \leq \infty$. Let \mathcal{L}_q be the collection of convex losses Lipschitz with respect to the ℓ_q -norm, and assume Θ contains the ℓ_p -ball of radius R_p . Then there is a numerical constant $c > 0$ such that for all $n \geq d$ and $G > 0$,*

$$\mathfrak{M}_n(\Theta, G \cdot \mathcal{L}_q) \geq c R_p G \frac{d^{1/2-1/p}}{\sqrt{n}}.$$

Proof We sketch the proof. The ℓ_p -ball of radius R_p contains the ℓ_∞ ball of radius $R_\infty := R_p/d^{1/p}$. The modification of the losses (17.3.8) in the proof of Corollary 17.3.11 guarantees that they are Lipschitz with respect to any ℓ_q -norm with appropriate constant, as the loss defined for 1-sparse $x \in \mathbb{R}^d$ by $\ell(\theta; x) = G|\theta_j - \text{sign}(x_j)R_p|$ when $x_j \neq 0$ is the unique non-zero element of x always satisfies $\|g\|_1 \leq G$ for $g \in \partial\ell(\theta; x)$. In particular, for $q^* = \frac{q}{q-1}$ conjugate to q , $\|g\|_{q^*} \leq \|g\|_1 \leq G$. Applying Corollary 17.3.11 gives

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq c \frac{R_p}{d^{1/p}} G \min \left\{ \frac{\sqrt{d}}{\sqrt{n}}, 1 \right\},$$

and when $n \geq d$, the first term in the minimum is smaller than the second. \square

Written differently by considering conjugates, if we let \mathcal{L} consist of losses ℓ with subgradients $g \in \partial\ell(\theta)$ implies $\|g\|_q \leq G_q$ for some $1 \leq q \leq 2$, and the domain $\Theta \supset \{\theta \in \mathbb{R}^d \mid \|\theta\|_p \leq R_p\}$, then

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq c R_p G_q \frac{d^{1/2-1/p}}{\sqrt{n}}$$

for $n \geq d$.

Comparing to the regret bounds for stochastic and online gradient methods (Corollary 17.2.2) shows that these results are sharp. Indeed, let $\Theta \subset \{\theta \in \mathbb{R}^d \mid \|\theta\|_2 \leq R_2\}$ and assume $\|g\|_2 \leq G_2$ for all subgradients g , and let θ_t be the iterates of the stochastic gradient method (17.2.2) with stepsize $\eta = \frac{R_2}{G_2} \frac{1}{\sqrt{n}}$, as in Corollary 17.2.2. Then the average $\hat{\theta}_n = \frac{1}{n} \sum_{t=1}^n \theta_t$ satisfies

$$\mathbb{E}[L_P(\hat{\theta}_n)] \leq L_P(\theta^*) + \frac{G_2 R_2}{\sqrt{n}} \quad \text{for any } \theta^* \in \Theta$$

as in the inequality (17.0.3). Notably, $\|g\|_2 \leq \|g\|_q$ for all $1 \leq q \leq 2$, and $\|\theta\|_2 \leq d^{1/2-1/p} \|\theta\|_p$ for $p \in [2, \infty]$. Thus the ℓ_p -ball of radius R_p is contained within the ℓ_2 -ball of radius $d^{1/2-1/p} R_p$, so the stochastic gradient method guarantees

$$\mathbb{E}[L_P(\hat{\theta}_n)] - L_P(\theta^*) \leq G_q R_p \cdot \frac{d^{1/2-1/p}}{\sqrt{n}}$$

whenever $\Theta \subset \{\theta \mid \|\theta\|_p \leq R_p\}$ and the losses satisfy $\|g\|_q \leq G_q$ for $g \in \partial\ell(\theta)$. Comparing to Corollary 17.3.12, these rates of convergence are evidently minimax optimal.

These results provide sharp lower bounds for certain cases and gradient/parameter space geometries, and in particular, the dual geometries when Θ is an ℓ_p -ball for some $p \geq 2$ and \mathcal{L} consists of functions Lipschitz with respect to the ℓ_p -norm. Here, we state (leaving the proof to the exercises) a result that naturally handles the geometry in which Θ is contained in an ℓ_p -ball for $p \leq 2$. We say a convex set $\Theta \subset \mathbb{R}^d$ is *orthosymmetric* if for $\theta \in \Theta$, $S\theta \in \Theta$ for any diagonal sign matrix S . Any ℓ_p -ball is orthosymmetric. For a vector of nonnegative parameters $\{G_j\}_{j=1}^d$, we let

$$\mathcal{L}(\{G_j\}) := \left\{ \text{convex } \ell : \mathbb{R}^d \rightarrow \mathbb{R} \text{ s.t. } g \in \partial \ell(\theta) \text{ satisfies } |g_j| \leq G_j \text{ for } j = 1, \dots, d \right\}, \quad (17.3.11)$$

that is, whose subgradients have coordinates bounded as specified. Any loss of the form $\ell(\theta; x) = \sum_{j=1}^d G_j |\theta_j - x_j|$ evidently belongs to the class (17.3.11), as do linear functions $\ell(\theta; x) = \langle x, \theta \rangle$ for vectors x satisfying $|x_j| \leq G_j$. The following theorem provides a lower bound for this class.

Theorem 17.3.13. *Let $\Theta \subset \mathbb{R}^d$ be an orthosymmetric convex set and \mathcal{L} be the class (17.3.11) of losses. Then*

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq \frac{1}{8\sqrt{n}} \sup_{\theta \in \Theta} \sum_{j=1}^d G_j |\theta_j|.$$

Exercise 17.11 steps through one approach to proving Theorem 17.3.13.

As a corollary to Theorem 17.3.13, we can provide minimax lower bounds for more general norms on Θ and the gradients. We say that $\|\cdot\|$ is orthosymmetric if $\|Sv\| = \|v\|$ for any vector v and diagonal matrix of signs S , which by inspection implies that the associated dual norm $\|\cdot\|_*$ is orthosymmetric. Then by carefully choosing the coordinate-Lipschitz parameters G_j in Theorem 17.3.13, we have the following corollary.

Corollary 17.3.14. *Let $\|\cdot\|$ be an orthosymmetric norm, and let $\Theta \subset \mathbb{R}^d$ an orthosymmetric convex set. Let \mathcal{L} consist of losses G -Lipschitz with respect to the norm $\|\cdot\|$, that is, for which $g \in \partial_\theta \ell(\theta)$ implies $\|g\|_* \leq G$. Then*

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq \frac{G}{8\sqrt{n}} \sup_{\theta \in \Theta} \|\theta\|.$$

That is, the norms $\|\cdot\|$ on Θ and the magnitude of the dual norms $\|\cdot\|_*$ of gradients necessarily appear in any lower bound.

When $\Theta = \{\theta \in \mathbb{R}^d \mid \|\theta\|_p \leq R_p\}$ is an ℓ_p -ball of radius R_p and \mathcal{L} consists of G_q -Lipschitz loss functions for the ℓ_p -norm, where $q = \frac{p}{p-1}$ is conjugate to p , then Corollary 17.3.14 implies the lower bound

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \gtrsim \frac{1}{\sqrt{n}} G_q R_p.$$

For $p > 1$, assuming that $\theta_1 = \mathbf{0}$ and we use a fixed stepsize $\eta > 0$, the p -norm algorithms as in Example 17.2.13 guarantee regret

$$\sum_{t=1}^n \ell_t(\theta_t) - \ell_t(\theta^*) \leq \frac{\|\theta^*\|_p^2}{2(p-1)\eta} + \frac{\eta}{2} n G_q^2 \leq \frac{R_p^2}{2(p-1)\eta} + \frac{\eta}{2} n G_q^2$$

for any sequence of losses $\ell_t \in \mathcal{L}$. Choosing $\eta = \frac{R_p}{G_q} \frac{1}{\sqrt{(p-1)n}}$ then yields the minimax upper bound

$$\mathbb{E} \left[L_P(\hat{\theta}_n) \right] - \inf_{\theta^* \in \Theta} L_P(\theta^*) \leq \frac{G_q R_p}{\sqrt{(p-1)n}},$$

where $\hat{\theta}_n = \frac{1}{n} \sum_{t=1}^n \theta_t$ is the average of the iterates of the mirror descent method 17.3 and we have applied Jensen's inequality as in (17.0.3). So, at least to the factor of $\sqrt{p-1}$, these methods achieve minimax optimal convergence. (We can in fact show that the factor $\sqrt{p-1}$ is necessary, though this is beyond our scope.)

17.3.4 A lower bound for high-dimensional stochastic optimization

We conclude our discussion of lower bounds for stochastic optimization by considering the particular geometry that arises in sparse problems, where Θ is an ℓ_1 ball and gradients g lie in ℓ_∞ balls. We have seen (inequalities (17.2.6) and Example 17.2.13) that $\sqrt{\log d/n}$ upper bounds the average regret and stochastic convergence rate. This, it turns out, is sharp.

Proposition 17.3.15. *Let $\Theta = \{\theta \mid \|\theta\|_1 \leq R_1\}$ be a scaled ℓ_1 -ball and \mathcal{L} consist of the losses satisfying $\|g\|_\infty \leq G_\infty$ for any $g \in \partial\ell(\theta)$. Then*

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq \frac{1}{24} G_\infty R_1 \min \left\{ \sqrt{\frac{\log(2d)}{n}}, 1 \right\}.$$

Proof For shorthand, let $G = G_\infty$ and $R = R_1$. Consider the sample space $\mathcal{X} = \{-1, 1\}^d$ and consider the linear losses

$$\ell(\theta; x) := G \langle \theta, x \rangle,$$

which evidently satisfy $\|\nabla \ell(\theta; x)\|_\infty = G \|x\|_\infty = G$. Define the packing set $\mathcal{V} := \{\pm e_j\}_{j=1}^d$ of the signed standard basis vectors and let P_v be the distribution on $X \in \{\pm 1\}^d$ with independent coordinates

$$X_j = \begin{cases} 1 & \text{w.p. } \frac{1+\delta v_j}{2} \\ -1 & \text{w.p. } \frac{1-\delta v_j}{2}. \end{cases}$$

We first demonstrate the separations we use to prove the lower bound, then apply Fano's method. Define $L_v(\theta) = \mathbb{E}_{P_v}[\ell(\theta; X)] = G \delta \langle v, \theta \rangle$, so that

$$\theta^v := \operatorname{argmin}_{\theta \in \Theta} L_v(\theta) = -Rv \quad \text{and} \quad L_v^* := \inf_{\theta \in \Theta} L_v(\theta) = -GR\delta.$$

For any two losses with $v \neq v'$, we see that

$$\inf_{\theta \in \Theta} \{L_v(\theta) + L_{v'}(\theta)\} = \inf_{\|\theta\|_1 \leq R} \{G\delta \langle v + v', \theta \rangle\} = -GR\delta \|v + v'\|_\infty \geq -GR\delta$$

as $\|v + v'\|_\infty = 0$ or 1 for $v \neq v' \in \{\pm e_j\}_{j=1}^d$. Thus, recalling the optimization distance (17.3.7), we have

$$d_{\text{opt}}(L_v, L_{v'}; \Theta) \geq -GR\delta + 2GR\delta = GR\delta.$$

Substituting this into the optimization-to-testing lower bound (17.3.6), we obtain that if we draw $V \sim \text{Uniform}(\mathcal{V})$ and, conditional on $V = v$, draw $X_1^n \stackrel{\text{iid}}{\sim} P_v$, then

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq \frac{GR\delta}{2} \inf_{\hat{v}} \mathbb{P}(\hat{v}(X_1^n) \neq V).$$

We may now apply Fano's method, which implies

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq \frac{GR\delta}{2} \left(1 - \frac{I(X_1^n; V) + \log 2}{\log(2d)} \right).$$

To bound the mutual information, consider the KL-divergence $D_{\text{kl}}(P_v \| P_{v'})$ for a pair $v \neq v'$. Let and $D_0(\delta) = D_{\text{kl}}(\text{Bernoulli}(\frac{1+\delta}{2}) \| \text{Bernoulli}(\frac{1}{2}))$ and $D_1(\delta) = D_{\text{kl}}(\text{Bernoulli}(\frac{1+\delta}{2}) \| \text{Bernoulli}(\frac{1-\delta}{2}))$. Then by inspection

$$D_{\text{kl}}(P_v \| P_{v'}) \leq \max \{2D_0(\delta), D_1(\delta)\} \leq D_1(\delta) \leq \frac{9}{4}\delta^2,$$

the final inequality valid for $0 \leq \delta \leq \frac{1}{2}$. So $I(X_1^n; V) \leq \frac{9n}{4}\delta^2$, and

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq \frac{GR\delta}{2} \left(1 - \frac{\log 2}{\log(2d)} - \frac{9n\delta^2}{4\log(2d)} \right).$$

Assuming that $d \geq 2$, we have $\log 2 / \log(2d) \leq \frac{1}{2}$, and so taking $\delta^2 = \frac{\log(2d)}{9n} \wedge \frac{1}{2}$ to make the testing-error component at least $\frac{1}{4}$ implies the lower bound

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq \frac{GR\delta}{8}.$$

Substituting the value of δ gives the proposition when $d \geq 2$.

In the case that $d = 1$, a straightforward argument via Le Cam's two-point method gives the claimed lower bound. \square

17.4 Online to batch conversions

The application (17.0.3) of Jensen's inequality shows that if an algorithm has a regret bound, then it enjoys a similar convergence rate for stochastic optimization. More precisely, if for a collection of losses \mathcal{L} a procedure has regret guarantee $\text{Reg}_n(\theta) \leq C_n$ for any sequence of losses $\ell_t \in \mathcal{L}$ and any $\theta \in \Theta$, then the average $\bar{\theta}_n := \frac{1}{n} \sum_{t=1}^n \theta_t$ satisfies $\mathbb{E}[L_P(\bar{\theta}_n)] \leq L_P(\theta) + C_n/n$ for any $\theta \in \Theta$. We can extend this to provide high-probability convergence guarantees as well using the martingale convergence theorems, especially the Azuma-Hoeffding inequality in Theorem 4.2.3.

In the next theorem, a prototypical result, we define the observed regret

$$\text{Reg}_n(\theta) = \sum_{i=1}^n \ell(\theta_i, Z_i) - \ell(\theta, Z_i)$$

for a sequence $Z_1^n \stackrel{\text{iid}}{\sim} P$. The theorem shows that this quantity provides a high-probability bound on the suboptimality gap on the population expected loss $L_P(\theta) := \mathbb{E}_P[\ell(\theta, Z)]$. The key is that, for any online optimization procedure, the i th point θ_i is a function only of the past random variables Z_1^{i-1} , so that Z_i provides fresh (independent) randomness at each iteration.

Theorem 17.4.1. *Let \mathcal{L} consist of convex losses, G -Lipschitz continuous with respect to the norm $\|\cdot\|$, and let Θ satisfy $\sup_{\theta, \theta' \in \Theta} \|\theta - \theta'\| \leq R$. Let θ_i be any sequence predicted by an online procedure. Then for any distribution P ,*

$$\mathbb{P} \left(L_P(\bar{\theta}_n) - L_P(\theta) \geq n^{-1} \text{Reg}_n(\theta) + 2\sqrt{2GRt} \right) \leq \exp(-nt^2).$$

Proof Let $\mathcal{F}_i := \{Z_1^i\}$ be the information available to the algorithm at time i , and let $\ell_i(\theta) = \ell(\theta; Z_i)$. We construct a sub-Gaussian martingale difference sequence using these (recall Definition 4.5), after which we may apply the Azuma-Hoeffding concentration bound (Theorem 4.2.3). Thus, define the martingale differences

$$D_i := (L_P(\theta_i) - L_P(\theta)) - (\ell_i(\theta_i) - \ell_i(\theta)),$$

which evidently are functions of \mathcal{F}_i and satisfy

$$\mathbb{E}[D_i \mid \mathcal{F}_{i-1}] = (L_P(\theta_i) - L_P(\theta)) - \mathbb{E}[\ell(\theta_i, Z_i) - \ell(\theta, Z_i) \mid Z_1^{i-1}] = 0,$$

because Z_i is independent of Z_1^{i-1} and θ_i is a function of Z_1^{i-1} . Additionally, under the Lipschitz and boundedness conditions on Θ and ℓ , we have $|D_i| \leq 2G \|\theta_i - \theta\| \leq 2GR$. So the D_i are thus a $4G^2R^2$ -sub-Gaussian martingale difference sequence, and Theorem 4.2.3 implies

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n D_i \geq t\right) \leq \exp\left(-\frac{nt^2}{8G^2R^2}\right)$$

for all $t \geq 0$. Now observe that

$$L_P(\bar{\theta}_n) - L_P(\theta) \leq \frac{1}{n} \sum_{i=1}^n (L_P(\theta_i) - L_P(\theta)) = \frac{1}{n} \sum_{i=1}^n D_i + \text{Reg}_n(\theta).$$

Evidently this implies the result. □

JCD Comment: Maybe add a bit of commentary around this and about why it is a good generalization bound. Compare to the uniform convergence guarantees in previous chapters.

17.5 More refined convergence guarantees

It is sometimes possible to give more refined bounds than those we have so far provided. As motivation, let us revisit Example 17.1.3, but suppose that one of the experts has no loss—that is, it makes perfect predictions. We might expect—accurately!—that we should attain better convergence guarantees using exponentiated weights, as the points w_t we maintain should quickly eliminate non-optimal experts.

To that end, we present a refined regret bound for the mirror descent algorithm 17.3 with the entropic regularization $\psi(w) = \sum_j w_j \log w_j$.

Proposition 17.5.1. *Let $\psi(w) = \sum_j w_j \log w_j$, and assume that the losses ℓ_t are such that their subgradients have all non-negative entries, that is, $g_t \in \partial \ell_t(w)$ implies $g_t \succeq 0$. For any such sequence of loss functions ℓ_t and any $w^* \in \Theta = \Delta_d$,*

$$\sum_{t=1}^n [\ell_t(w_t) - \ell_t(w^*)] \leq \frac{\log d}{\eta} + \frac{\eta}{2} \sum_{t=1}^n \sum_{j=1}^d w_{t,j} g_{t,j}^2.$$

While as stated, the bound of the proposition does not look substantially more powerful than Corollary 17.2.12, but a few remarks will exhibit its consequences. We prove the proposition in Section 17.5.1 to come.

First, we note that because $w_t \in \Delta_d$, we will *always* have $\sum_j w_{t,j} g_{t,j}^2 \leq \|g_t\|_\infty^2$. So certainly the bound of Proposition 17.5.1 is never worse than that of Corollary 17.2.12. Sometimes this can be made tighter, however, as exhibited by the next corollary, which applies (for example) to the experts setting of Example 17.1.3. More specifically, we have d experts, each suffering losses in $[0, 1]$, and we seek to predict with the best of the d experts.

Corollary 17.5.2. *Consider the linear online convex optimization setting, that is, where $\ell_t(w_t) = \langle g_t, w_t \rangle$ for vectors g_t , and assume that $g_t \in \mathbb{R}_+^d$ with $\|g_t\|_\infty \leq 1$. In addition, assume that we know an upper bound L_n^* on $\sum_{t=1}^n \ell_t(w^*)$. Then taking the stepsize $\eta = \min\{1, \sqrt{\log d}/\sqrt{L_n^*}\}$, we have*

$$\sum_{t=1}^n [\ell_t(w_t) - \ell_t(w^*)] \leq 3 \max \left\{ \log d, \sqrt{L_n^* \log d} \right\}.$$

Note that when $\ell_t(w^*) = 0$ for all w^* , which corresponds to a perfect expert in Example 17.1.3, the upper bound becomes constant in n , yielding $3 \log d$ as a bound on the regret. Unfortunately, in our bound of Corollary 17.5.2, we had to assume that we *knew* ahead of time a bound on the loss of the best predictor w^* , which is unrealistic in practice. There are a number of techniques for dealing with such issues, including a standard one in the online learning literature known as the *doubling* trick. We explore some in the exercises.

Proof First, we note that $\sum_j w_j g_{t,j}^2 \leq \langle w, g_t \rangle$ for any nonnegative vector w , as $g_{t,j} \in [0, 1]$. Thus, Proposition 17.5.1 gives

$$\sum_{t=1}^n [\ell_t(w_t) - \ell_t(w^*)] \leq \frac{\log d}{\eta} + \frac{\eta}{2} \sum_{t=1}^n \langle w_t, g_t \rangle = \frac{\log d}{\eta} + \frac{\eta}{2} \sum_{t=1}^n \ell_t(w_t).$$

Rearranging via an algebraic manipulation, this is equivalent to

$$\left(1 - \frac{\eta}{2}\right) \sum_{t=1}^n [\ell_t(w_t) - \ell_t(w^*)] \leq \frac{\log d}{\eta} + \frac{\eta}{2} \sum_{t=1}^n \ell_t(w^*).$$

Take $\eta = \min\{1, \sqrt{\log d/L_n^*}\}$. Then if $\sqrt{\log d/L_n^*} \leq 1$, we have that the right hand side of the above inequality becomes $\sqrt{L_n^* \log d} + \frac{1}{2} \sqrt{L_n^* \log d}$. On the other hand, if $L_n^* < \log d$, then the right hand side of the inequality becomes $\log d + \frac{1}{2} L_n^* \leq \frac{3}{2} \log d$. In either case, we obtain the desired result by noting that $1 - \frac{\eta}{2} \geq \frac{1}{2}$. \square

17.5.1 Proof of Proposition 17.5.1

Our proof relies on a technical lemma, after which the derivation is a straightforward consequence of Lemma 17.2.16. We first state the technical lemma, which applies to the update that the exponentiated gradient procedure makes.

Lemma 17.5.3. *Let $\psi(x) = \sum_j x_j \log x_j$, and let $x, y \in \Delta_d$ be defined by*

$$y_i = \frac{x_i \exp(-\eta g_i)}{\sum_j x_j \exp(-\eta g_j)},$$

where $g \in \mathbb{R}_+^d$ is non-negative. Then

$$-\frac{1}{\eta}D_\psi(y, x) + \langle g, x - y \rangle \leq \frac{\eta}{2} \sum_{i=1}^d g_i^2 x_i.$$

Deferring the proof of the lemma, we note that it precisely applies to the setting of Lemma 17.2.16. Indeed, with a fixed stepsize η , we have

$$\sum_{t=1}^n \ell_t(w_t) - \ell_t(w^*) \leq \frac{1}{\eta}D_\psi(w^*, w_1) + \sum_{t=1}^n \left[-\frac{1}{\eta}D_\psi(w_{t+1}, w_t) + \langle g_t, w_t - w_{t+1} \rangle \right].$$

Earlier, we used the strong convexity of ψ to eliminate the gradient terms $\langle g_t, w_t - w_{t+1} \rangle$ using the bregman divergence D_ψ . This time, we use Lemma 17.2.16: setting $y = w_{t+1}$ and $x = w_t$ yields the bound

$$\sum_{t=1}^n \ell_t(w_t) - \ell_t(w^*) \leq \frac{1}{\eta}D_\psi(w^*, w_1) + \sum_{t=1}^n \frac{\eta}{2} \sum_{i=1}^d g_{t,i}^2 w_{t,i}$$

as desired.

Proof of Lemma 17.5.3 We begin by noting that a direct calculation yields $D_\psi(y, x) = D_{\text{kl}}(y \| x) = \sum_i y_i \log \frac{y_i}{x_i}$. Substituting the values for x and y into this expression, we have

$$\sum_i y_i \log \frac{y_i}{x_i} = \sum_i y_i \log \left(\frac{x_i \exp(-\eta g_i)}{x_i (\sum_j \exp(-\eta g_j) x_j)} \right) = -\eta \langle g, y \rangle - \sum_i y_i \log \left(\sum_j x_j e^{-\eta g_j} \right).$$

Now we use a Taylor expansion of the function $g \mapsto \log(\sum_j x_j e^{-\eta g_j})$ around the point 0. If we define the vector $p(g)$ by $p_i(g) = x_i e^{-\eta g_i} / (\sum_j x_j e^{-\eta g_j})$, then

$$\log \left(\sum_j x_j e^{-\eta g_j} \right) = \log(\langle \mathbf{1}, x \rangle) - \eta \langle p(0), g \rangle + \frac{\eta^2}{2} g^\top (\text{diag}(p(\tilde{g})) - p(\tilde{g})p(\tilde{g})^\top) g,$$

where $\tilde{g} = \lambda g$ for some $\lambda \in [0, 1]$. Noting that $p(0) = x$ and $\langle \mathbf{1}, x \rangle = \langle \mathbf{1}, y \rangle = 1$, we obtain

$$D_\psi(y, x) = -\eta \langle g, y \rangle + \log(1) + \eta \langle g, x \rangle - \frac{\eta^2}{2} g^\top (\text{diag}(p(\tilde{g})) - p(\tilde{g})p(\tilde{g})^\top) g,$$

whence

$$-\frac{1}{\eta}D_\psi(y, x) + \langle g, x - y \rangle \leq \frac{\eta}{2} \sum_{i=1}^d g_i^2 p_i(\tilde{g}). \quad (17.5.1)$$

Lastly, we claim that the function

$$s(\lambda) = \sum_{i=1}^d g_i^2 \frac{x_i e^{-\lambda g_i}}{\sum_j x_j e^{-\lambda g_j}}$$

is non-increasing on $\lambda \in [0, 1]$. Indeed, we have

$$s'(\lambda) = \frac{(\sum_i g_i x_i e^{-\lambda g_i})(\sum_i g_i^2 x_i e^{-\lambda g_i})}{(\sum_i x_i e^{-\lambda g_i})^2} - \frac{\sum_i g_i^3 x_i e^{-\lambda g_i}}{\sum_i x_i e^{-\lambda g_i}} = \frac{\sum_{ij} g_i g_j^2 x_i x_j e^{-\lambda g_i - \lambda g_j} - \sum_{ij} g_i^3 x_i x_j e^{-\lambda g_i - \lambda g_j}}{(\sum_i x_i e^{-\lambda g_i})^2}.$$

Using the Fenchel-Young inequality, we have $ab \leq \frac{1}{3}|a|^3 + \frac{2}{3}|b|^{3/2}$ for any a, b , so $g_i g_j^2 \leq \frac{1}{3}g_i^3 + \frac{2}{3}g_j^3$. This implies that the numerator in our expression for $s'(\lambda)$ is non-positive. Thus we have $s(\lambda) \leq s(0) = \sum_{i=1}^d g_i^2 x_i$, which gives the result when combined with inequality (17.5.1). \square

17.6 Exercises

Exercise 17.1: We consider the doubling trick, a frequently used technique in online learning to allow good performance of online learning procedures even without knowledge of the number of steps n they will be run. For this question, we define the regret in the usual way as

$$\text{Reg}_n := \sup_{w^* \in \Theta} \sum_{t=1}^n [\ell_t(w_t) - \ell_t(w^*)].$$

- (a) Suppose that we have a procedure (algorithm) $A(\eta)$ parameterized by the real value $\eta \geq 0$ (usually, this is simply a stepsize) that achieves the regret bound

$$\text{Reg}_n \leq \frac{r^2}{2\eta} + \frac{\eta}{2} L^2 n$$

where r and L are known constants. Consider the following procedure, which proceeds in epochs $k = 1, 2, \dots$, each of which lasts for $n_k = 2^k$ steps. At the start of epoch k , restart the algorithm $A(\eta)$ with parameter choice $\eta_k = \frac{r}{L\sqrt{2^k}}$, and run the algorithm with this choice of parameter for 2^k steps. Show that

$$\text{Reg}_n \leq C \cdot Lr\sqrt{n},$$

where C is some numerical constant (in our solution, we have $C \leq 2/(\sqrt{2} - 1)$).

- (b) Now we consider a slightly more restrictive setting, but we obtain better guarantees. Consider the mixture of experts problem, in d experts suffer losses in $[0, 1]$ at each timestep; we let $g_t \in [0, 1]^d$ denote the loss vector. We play a mixture of experts $w_t \in \Theta = \Delta_d$, suffering (expected) loss $\ell_t(w_t) = \langle g_t, w_t \rangle$. In the course notes, we show that in this setting

$$\text{Reg}_n = \sum_{t=1}^n [\ell_t(w_t) - \ell_t(w^*)] \leq \frac{\log d}{\eta} + \frac{\eta}{2} \sum_{t=1}^n \sum_{j=1}^d w_{t,j} g_{t,j}^2$$

when using the exponential weights algorithm with stepsize η , which in turn implies

$$\sum_{t=1}^n [\ell_t(w_t) - \ell_t(w^*)] \leq \left(1 - \frac{\eta}{2}\right)^{-1} \left[\frac{\log d}{\eta} + \frac{\eta}{2} \sum_{t=1}^n \ell_t(w^*) \right]$$

for any $w^* \in \Theta$. Consider the following procedure, which proceeds in epochs $k = 1, 2, \dots$, within each of which we perform exponential weights with stepsize $\eta_k = \min\{1, \sqrt{\log d}/2^k\}$. Let E_k denote those times t belonging to epoch k , which correspond to times when we run exponential weights with parameter η_k . Define $L^{(k)} = \min_j \sum_{t \in E_k} g_{t,j}$ to be the loss incurred by the best expert in epoch k as the procedure runs, and continue epoch k until the best expert's loss in epoch k satisfies $L^{(k)} \geq 4^k$. Then begin a new epoch. Show that with this procedure,

$$\sum_{t=1}^n [\ell_t(w_t) - \ell_t(w^*)] \leq C_1 \log \log d \cdot \log d + C_2 \sqrt{\log d \cdot \sum_{t=1}^n \ell_t(w^*)}$$

for numerical constants C_1 and C_2 (we obtain $C_1 \leq 3$ and $C_2 \leq 8\sqrt{2}$).

Exercise 17.2 (Min-max games and regret): The saddle point problem, or min-max game problem, considers solving

$$\underset{x \in X}{\text{minimize}} \sup_{y \in Y} L(x, y), \quad (17.6.1)$$

where L is convex in its first argument and concave in its second, and X, Y are convex sets. (See Appendix C.4 for a general treatment of these problems.) We say a point $(x^*, y^*) \in X \times Y$ is a saddle point if

$$\sup_{y \in Y} L(x^*, y) \leq L(x^*, y^*) \leq \inf_{x \in X} L(x, y^*).$$

(a) Show that if a saddle point exists, then

$$\inf_{x \in X} \sup_{y \in Y} L(x, y) = \sup_{y \in Y} \inf_{x \in X} L(x, y) = L(x^*, y^*).$$

Now we show how online convex optimization can prove that saddle points exist for some problems. Assume that X and Y are compact and convex, and that for each $x_0 \in X$, $L(x_0, y)$ is 1-Lipschitz and concave in y , and for each $y_0 \in Y$, $L(x, y_0)$ is 1-Lipschitz and convex and subdifferentiable in x . Consider the following online game: at iteration t , the x player chooses $x_t \in X$, and then the Y player chooses the “best response”

$$y_t \in \operatorname{argmax}_{y \in Y} L(x_t, y).$$

The goal of the x player is to achieve low regret with respect to any fixed $x^* \in X$.

(b) Define $f_t(x) = L(x, y_t)$. Give a strategy for the x player so that for any T and any $x^* \in X$,

$$\sum_{t=1}^T f_t(x_t) - f_t(x^*) \leq O(1)\sqrt{T}.$$

(c) Show that $\bar{x}_T = \frac{1}{T} \sum_{t=1}^T x_t$ and $\bar{y}_T = \frac{1}{T} \sum_{t=1}^T y_t$ satisfy

$$\sup_{y \in Y} L(\bar{x}_T, y) \leq L(x^*, \bar{y}_T) + O(1)/\sqrt{T}$$

for any $x^* \in X$.

(d) Show that there exists a saddle point for the min-max problem (17.6.1) when L is Lipschitz and X and Y are compact, as above.

Exercise 17.3 (von Neumann’s minimax theorem for zero-sum games): Let $A \in \mathbb{R}^{m \times n}$ be an arbitrary matrix and X and Y be compact convex sets. Show that

$$\inf_{x \in X} \sup_{y \in Y} \langle x, Ay \rangle = \sup_{y \in Y} \inf_{x \in X} \langle x, Ay \rangle,$$

and that there exists a saddle point $(x^*, y^*) \in X \times Y$ for which

$$\sup_{y \in Y} \langle x, Ay \rangle \leq \langle x^*, Ay^* \rangle \leq \inf_{x \in X} \langle x, Ay^* \rangle.$$

In the language of games, there are strategies (x^*, y^*) for which it is unimportant which player plays first.

Exercise 17.4 (Second-order strong convexity conditions): Let f be twice continuously differentiable on its domain, which you may assume to be open.

- (a) Assume that f is strongly convex with respect to the norm $\|\cdot\|$. Show that for all $x \in \text{dom } f$, $u^\top \nabla^2 f(x) u \geq \|u\|^2$ for all u .
- (b) Show that if $u^\top \nabla^2 f(x) u \geq \|u\|^2$ for all u and $x \in \text{dom } f$, then f is strongly convex with respect to the norm $\|\cdot\|$.

Hint. See Proposition C.1.5 in Appendix C.1.

Exercise 17.5 (The strong convexity of p -norms): In this question, we use the results of Exercise 17.4 to show that for $1 < p < 2$, $f_p(w) := \frac{1}{2(p-1)} \|w\|_p^2$ is strongly convex with respect to the norm $\|\cdot\|_p$.

- (a) Define $\psi(t) = \frac{1}{2(p-1)} t^{2/p}$ and $\phi(t) = |t|^p$. Let $H(w) = \nabla^2 f_p(w)$. Show that

$$H_{ii}(w) = \psi'' \left(\|w\|_p^p \right) (\phi'(w_i))^2 + \psi' \left(\|w\|_p^p \right) \phi''(w_i)$$

and for $i \neq j$,

$$H_{ij}(w) = \psi'' \left(\|w\|_p^p \right) \phi'(w_i) \phi'(w_j).$$

- (b) Show that for all u ,

$$u^\top \nabla^2 f_p(w) u \geq \|w\|_p^{2-p} \sum_{j=1}^d \frac{u_j^2}{|w_j|^{2-p}}.$$

- (c) Use Hölder's inequality to demonstrate that $u^\top \nabla^2 f_p(w) u \geq \|u\|_p^2$ for all u .

Exercise 17.6: Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex.

- (a) Show that if f is G -Lipschitz with respect to the norm $\|\cdot\|$, meaning $|f(x) - f(y)| \leq G \|x - y\|$, then for all $g \in \partial f(x)$ and all x , $\|g\|_* \leq G$.
- (b) Show that if for all $g \in \partial f(x)$ and all x , $\|g\|_* \leq G$, then f is G -Lipschitz with respect to the norm $\|\cdot\|$.

Hint. Use Corollary B.3.19 and Proposition B.3.20 in Appendix B.3.4.

Exercise 17.7 (The naming of mirror descent): Let $\psi : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be a Legendre distance generating function, meaning that it is strictly convex, continuously differentiable, and satisfies the conditions (14.3.3). Show that for any $g \in \mathbb{R}^d$, $\eta > 0$, and $\theta_0 \in \text{dom } \psi$,

$$\theta_\eta = \underset{\theta}{\operatorname{argmin}} \left\{ \langle g, \theta \rangle + \frac{1}{\eta} D_\psi(\theta, \theta_0) \right\}.$$

satisfies

$$\theta_\eta = \nabla \psi^* (\nabla \psi(\theta_0) - \eta g).$$

We therefore may think of the update (17.2.5) when we use the first-order model $\hat{\ell}(\theta) = \ell(\theta_0) + \langle g, \theta - \theta_0 \rangle$ as performing a type of “mirror” operation: we transform θ_0 into a dual version $\nabla \psi(\theta_0)$, where gradients belong to the dual space, update the parameter, then “reflect” back into the primal space via $\nabla \psi^* : \mathbb{R}^d \rightarrow \text{dom } \psi$. *Hint.* Use Corollary 14.3.1.

Exercise 17.8: In this question, we show how to compute the p -norm-based mirror descent update (17.2.5).

- (a) Let $h(x) = \frac{1}{2\eta} \|x\|^2$. Show that the conjugate $h^*(y) = \frac{\eta}{2} \|y\|_*^2$.
- (b) Let $\psi(\theta) = \frac{1}{2} \|\theta\|_p^2$, where $1 < p \leq 2$. Let $\theta_0 \in \mathbb{R}^d$, and for a fixed $g \in \mathbb{R}^d$ define

$$\theta_\eta = \operatorname{argmin}_{\theta} \left\{ \langle g, \theta \rangle + \frac{1}{\eta} D_\psi(\theta, \theta_0) \right\}.$$

Show that if we let w have entries

$$w_j = \frac{1}{\|\theta_0\|^{p-2}} |\theta_{0,j}|^{p-1} \operatorname{sign}(\theta_{0,j}) - \eta g_j$$

then for $q = \frac{p}{p-1}$ conjugate to p ,

$$\theta_\eta = \frac{1}{\|w\|_q^{q-2}} [|w_j|^{q-1} \operatorname{sign}(w_j)]_{j=1}^d.$$

Hint. Use Exercise 17.7.

Exercise 17.9: Let $\psi(u) = \frac{1}{2} \|u\|_p^2$, where $1 < p < \infty$.

- (a) Show that for any u, v ,

$$D_\psi(u, v) \leq \frac{1}{2} \|u\|_p^2 + \frac{1}{2} \|u - v\|_p^2.$$

- (b) Let $\Theta \subset \mathbb{R}^d$ satisfy that $\|\theta - \theta_0\|_p \leq R_p$ for all $\theta \in \Theta$. Show that if $\psi(u) = \frac{1}{2} \|u - \theta_0\|_p^2$, then

$$\sup_{\theta, \theta' \in \Theta} D_\psi(\theta, \theta') \leq \frac{5}{2} R_p^2.$$

Exercise 17.10 (On adaptive stepsizes): Let ψ be a distance generating function strongly convex with respect to a norm $\|\cdot\|$ over Θ , and let θ_t be generated by the iteration (17.2.5). Fix $\eta > 0$, and define the t th stepsize

$$\eta_t = \frac{\eta}{\sqrt{\sum_{\tau=1}^t \|g_\tau\|_*^2}}, \quad (17.6.2)$$

which is computable at iteration t . Assume that $D_\psi(\theta^*, \theta) \leq R^2$ for all $\theta \in \Theta$.

- (a) Show that

$$\sum_{t=1}^n \ell_t(\theta_t) - \ell_t(\theta^*) \leq \frac{1}{\eta} R^2 \sqrt{\sum_{t=1}^n \|g_t\|_*^2} + \eta \sqrt{\sum_{t=1}^n \|g_t\|_*^2}.$$

- (b) Show that with the fixed stepsize multiplier $\eta = R$,

$$\sum_{t=1}^n \ell_t(\theta_t) - \ell_t(\theta^*) \leq \sqrt{2} \inf_{\alpha} \left\{ \frac{R^2}{\alpha} + \frac{\alpha}{2} \sum_{t=1}^n \|g_t\|_*^2 \right\},$$

so that the adaptive steps (17.6.2) enjoy a type of post-hoc optimality guarantee.

Exercise 17.11 (Coordinate-based lower bounds [73]): This question explores using Assouad's method for lower bounds in stochastic convex optimization, providing a proof of Theorem 17.3.13. Recall that a convex set $\Theta \subset \mathbb{R}^d$ is orthosymmetric if for any $\theta \in \Theta$, $S\theta \in \Theta$ for any diagonal sign matrix S . Let $G_j \geq 0$ and \mathcal{L} be the coordinate-bounded loss class (17.3.11).

- (a) Fix an arbitrary vector $a \in \mathbb{R}^d$. Define the loss $\ell : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$ by

$$\ell(\theta; x) := \sum_{j=1}^d G_j |\theta_j - a_j x_j|.$$

Show that for each $x \in \mathbb{R}^d$, $\ell(\cdot, x)$ belongs to the class \mathcal{L} .

- (b) Use a variant of Assouad's method to prove Theorem 17.3.13. That is, show that there exists a numerical constant $c > 0$ such that for any orthosymmetric convex parameter space Θ ,

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq c \frac{1}{\sqrt{n}} \sup_{\theta \in \Theta} \sum_{j=1}^d G_j |\theta_j|.$$

(It is possible but perhaps tedious to obtain a constant $c \geq \frac{1}{8}$.) *Hint.* Define distributions P_v on $X \in \{-1, 1\}^d$, indexed by $v \in \{-1, 1\}^d$, so that $X \sim P_v$ has independent coordinates $X_j \sim \text{Bernoulli}(\frac{1+v_j\delta}{2})$. Take $a \in \Theta$ and use the losses in the previous part.

Exercise 17.12 (A geometric lower bound): Prove Corollary 17.3.14.

Exercise 17.13 (Sharper constants in Proposition 17.3.15): In this exercise, we trace the argument of Proposition 17.3.15 to obtain sharper constants as $d \uparrow \infty$.

- (a) Let $D_0(\delta) = D_{\text{kl}}(\text{Bernoulli}(\frac{1+\delta}{2}) \parallel \text{Bernoulli}(\frac{1}{2}))$ and $D_1(\delta) = D_{\text{kl}}(\text{Bernoulli}(\frac{1+\delta}{2}) \parallel \text{Bernoulli}(\frac{1-\delta}{2}))$. For the construction of P_v in the proof of Proposition 17.3.15, show that

$$I(X_1^n; V) \leq n \left[\frac{d-1}{d} 2D_0(\delta) + \frac{1}{d} D_1(\delta) \right] \leq n \left[\left(1 - \frac{1}{d}\right) \frac{9}{8} \delta^2 + \frac{9}{4d} \delta^2 \right],$$

the second inequality valid for $0 \leq \delta \leq \frac{1}{2}$.

- (b) Conclude that as $d \uparrow \infty$,

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq \frac{G_\infty R_1 \delta}{2} \left(1 - o(1) - \frac{9n\delta^2}{8 \log(2d)} \right).$$

- (c) Show that for n and d scaling so that $\log(2d)/n \rightarrow 0$ but $d \rightarrow \infty$ as $n \rightarrow \infty$,

$$\mathfrak{M}_n(\Theta, \mathcal{L}) \geq \frac{4(1-o(1))}{9\sqrt{6}} \cdot G_\infty R_1 \sqrt{\frac{\log(2d)}{n}} > \frac{11}{61} \cdot G_\infty R_1 \sqrt{\frac{\log(2d)}{n}}.$$

JCD Comment: Next question should go with banditos

Exercise 17.14 (An empirical comparison of Bandit algorithms): In this question, you will investigate three algorithms for solving the Bandit problem: the Upper Confidence Bound algorithm

(UCB), Thompson sampling (also known as Posterior Sampling), and exponential gradient. You will attempt to maximize the reward achieved by the algorithms (note that in the notes, we sometimes maximize and sometimes minimize; make sure you have your signs correct!).

In particular, set the rewards for the arms in the following way:

- i. Let $\theta_1 = \frac{1}{2}$ and $\theta_2 = \frac{1}{2} - \epsilon, \dots, \theta_K = \frac{1}{2} - \epsilon$.
- ii. When arm j is sampled, return $Y = 1$ with probability θ_j and $Y = 0$ with probability $1 - \theta_j$.

Now, repeat the following experiment with the values

- (a) $K = 10$, $\epsilon = .1$, and $n = 10^6$ steps
- (b) $K = 10$, $\epsilon = .02$, and $n = 10^6$ steps.

Perform Thompson sampling (Example 18.3.4 in the notes) assuming that the prior on θ is to have each coordinate independent with **Beta**(1, 1) distribution. Perform UCB with the confidence parameter $\delta_t = 1/\sqrt{t}$ (Algorithm 18.1 in the notes) and the appropriate choice of the sub-Gaussian parameter σ^2 (*Hint*: use Hoeffding's lemma for σ^2). Perform exponentiated gradient (Algorithm 18.6) using the optimal stepsize choice η , when assuming that $\sigma^2 = \frac{1}{2}$ in the bound.

Plot your results for each of experiments (a) and (b). Which algorithm do you prefer?

Chapter 18

Exploration, exploitation, and bandit problems

Consider the following problem: we have a possible treatment for a population with a disease, but we do not know whether the treatment will have a positive effect or not. We wish to evaluate the treatment to decide whether it is better to apply it or not, and we wish to optimally allocate our resources to attain the best outcome possible. There are challenges here, however, because for each patient, we may only observe the patient's behavior and disease status in one of two possible states—under treatment or under control—and we wish to allocate as few patients to the group with worse outcomes (be they control or treatment) as possible. This balancing act between exploration—observing the effects of treatment or non-treatment—and exploitation—giving treatment or not as we decide which has better palliative outcomes—underpins the challenges in this chapter.

Our main focus will be variants of the K -armed bandit problem, so named because we imagine a player in a casino, choosing between K different slot machines. As this is a casino, the player will surely lose eventually, hence the bandit moniker. Each machine has a different and unknown reward distribution. The player wishes to put as much money as possible into the machine with the greatest expected reward. There is a substantial literature in statistics, operations research, economics, game theory, and computer science on variants of the problems we consider here.

18.1 The multi-armed bandit problem

In the most basic setting—and we shall elaborate this later—We consider the following sequential decision making scenario. We assume that there are K distributions P_1, \dots, P_K on \mathbb{R} , which we identify with K random variables $Y(1), \dots, Y(K)$. In this basic setting, each random variable $Y(i)$ has mean μ_i and is σ^2 -sub-Gaussian, meaning that

$$\mathbb{E} [\exp (\lambda(Y(i) - \mu_i))] \leq \exp \left(\frac{\lambda^2 \sigma^2}{2} \right). \quad (18.1.1)$$

The goal is to find the index i with the maximal mean μ_i without sampling sub-optimal distributions P_i (or random variables $Y(i)$) too often.

We consider this in an online setting, proceeding iteratively for $t = 1, 2, \dots$, where at each iteration t of the process, the player takes an action $A_t \in \{1, \dots, K\}$, then, conditional on $i = A_t$, observes a reward $Y_t(i)$ drawn independently from the distribution P_i . The key is that given the past, the action A_t is independent of the vector $Y_t := (Y_t(1), \dots, Y_t(K))$, and at time t , we can

never observe the entire feedback vector. Then the goal is to minimize the realized regret after n steps, which is

$$\text{Reg}_n := \sum_{t=1}^n \mu_{i^*} - \mu_{A_t}, \quad (18.1.2)$$

where $i^* \in \text{argmax}_i \mu_i$ so $\mu_{i^*} = \max_i \mu_i$. While the regret (18.1.2) involves only the means of random variables, it is still random, so we generally seek to give bounds on its expectation or high-probability guarantees on its value. In this chapter, we generally focus for simplicity on the expected regret,

$$\text{Reg}_n := \mathbb{E} \left[\sum_{t=1}^n \mu_{i^*} - \mu_{A_t} \right], \quad (18.1.3)$$

where the expectation is taken over any randomness in the player's actions A_t and in the repeated observations of the random variables $Y(1), \dots, Y(K)$.

Example 18.1.1 (Potential outcomes): Consider estimating the effect of a particular treatment on some disease. In a simple version of this setting, the set of actions is $\{0, 1\}$, corresponding to treatment or control. In a particular model for estimation of causal effects, the eponymous *Neyman-Rubin causal model*, we imagine an individual i as having *potential outcomes* $Y_i := (Y_i(0), Y_i(1))$, where $Y_i(0)$ indicates the individual's response under the control (no treatment), while $Y_i(1)$ indicates the individual's response to treatment.

These outcomes are *potential* because we can never view both: we may only observe one or the other, as a patient cannot be both treatment and control. When the action A_i (assignment to treatment or control) for patient i is independent of any particular characteristics of the patient, then

$$\mathbb{E}[Y_i(a) \mid A_i = a] = \mathbb{E}[Y_i(a)] =: \mu_a \quad (18.1.4)$$

is the expected response for the patient for treatment ($a = 1$) or control ($a = 0$), and $\mathbb{E}[Y_i(1) - Y_i(0)]$ is the (unobservable!) treatment effect for the patient. (Here, we elide a small detail: we are thinking of the particular individual i as a random representative individual, not attempting to estimate anything particular about them.) If we choose A_i in a way that depends on the individual i , then we may have confounding: $\mathbb{E}[Y_i(a) \mid A_i = a] \neq \mathbb{E}[Y_i(a)]$.

Given a group of n individuals, a fully randomized trial to estimate the treatment effect chooses $n/2$ of the individuals to receive treatment and $n/2$ to be in the control, uniformly at random. Then

$$\hat{\tau} := \frac{1}{n/2} \sum_{i:A_i=1} Y_i(A_i) - \frac{1}{n/2} \sum_{i:A_i=0} Y_i(A_i)$$

is an unbiased estimator for $\tau = \mathbb{E}[Y(1) - Y(0)]$, with $\mathbb{P}(|\hat{\tau} - \tau| \geq t) \leq \exp(-c \frac{nt^2}{\sigma^2})$ in the sub-gaussian setting (18.1.1). This treatment assignment strategy, while effective for estimating the treatment effect, typically incurs very high regret (18.1.3), at least if the treatment is effective or has strong negative outcomes, because it allocates too many individuals to either the control or treatment arm of the study. \diamond

JCD Comment: Connect this with Example 12.2.5

Example 18.1.1 highlights many of the central aspects of bandit-like problems. First, they have deep connections with causal reasoning: by selecting an action (or arm) A , we are intervening in a system, with a goal to identify the best intervention. Second, the online nature can be critical:

in medical settings, for example, it would unethical to continue a drug trial if it were ineffective. Finally, the example implicitly shows one of the major benefits of the online scenario: because we *select* the action A_t without any information about the outcomes $Y_t = (Y_t(1), \dots, Y_t(K))$, we can avoid issues of confounding that would arise trying to estimate or learn from offline data.

18.2 Confidence-based algorithms

A natural first strategy to consider is one based on confidence intervals with slight optimism. Roughly, if we believe the true mean μ_i for an arm i lies within $[\hat{\mu}_i - c_i, \hat{\mu}_i + c_i]$, where c_i is some interval (whose length decreases with time t), then we optimistically “believe” that the value of arm i is $\hat{\mu}_i + c_i$; then at iteration t , as our action A_t we choose the arm whose optimistic mean is the highest, thus hoping to maximize our received reward.

This strategy lies at the heart of the Upper Confidence Bound (UCB) family of algorithms [12], a simple variant of which we describe here. Before continuing, we recall the standard result on sub-Gaussian random variables of Corollary 4.1.10 in our context, though we require a somewhat more careful calculation because of the sequential nature of our process. Let

$$N_t(i) = \text{card}\{\tau \leq t \mid A_\tau = i\}$$

denote the number of times that arm i has been pulled by time t of the bandit process, and define

$$\hat{\mu}_t(i) := \frac{1}{N_t(i)} \sum_{\tau \leq t, A_\tau = i} Y_\tau(i),$$

to be the running average of the rewards of arm i at time t (computed only on those instances in which arm i was selected). The sequential nature of the bandit process coupled with the (conditionally) independent randomness in sample $Y_t(i) \stackrel{\text{iid}}{\sim} P_i$ when $A_t = i$ in our version of the bandit game implies an important distributional equality, which underpins our analyses throughout:

Lemma 18.2.1. *Let $Y'_\tau(i) \stackrel{\text{iid}}{\sim} P_i$, $\tau = 1, 2, \dots$, be independent copies of the random variables $Y_\tau(i)$. Then $Y'_\tau(i)$ is independent of $N_t(i)$ for all t and τ , and*

$$(\hat{\mu}_t(i), N_t(i)) \stackrel{\text{dist}}{=} (\hat{\mu}'_t(i), N_t(i)), \quad (18.2.1)$$

where $\hat{\mu}'_t(i) = \frac{1}{N_t(i)} \sum_{\tau: A_\tau = i} Y'_\tau(i)$ is the empirical mean of the copies $Y'_\tau(i)$ for those steps when arm i is selected.

We prove the claim (18.2.1) in the Appendix 18.6.1 to this chapter.

The distributional equality (18.2.1) implies that the means $\hat{\mu}_t(i)$ are sub-Gaussian, even conditional on the number of pulls $N_t(i)$. Indeed, we observe that for $\lambda \in \mathbb{R}$ and any $m \in \mathbb{N}$,

$$\mathbb{E}[\exp(\lambda m \hat{\mu}_t(i)) \mid N_t(i) = m] \stackrel{(\star)}{=} \mathbb{E}[\exp(\lambda m \hat{\mu}'_t(i)) \mid N_t(i) = m] \leq \exp\left(\frac{m \lambda^2 \sigma^2}{2}\right), \quad (18.2.2)$$

where equality (\star) follows from Lemma 18.2.1, while the inequality follows because conditional on the event $N_t(i) = m$, we have

$$\hat{\mu}'_t(i) \stackrel{\text{dist}}{=} \frac{1}{m} \sum_{t=1}^m Y'_t(i) \quad \text{for } Y'_t(i) \stackrel{\text{iid}}{\sim} P_i.$$

Thus, by taking conditional expectations over the value $N_t(i)$, we immediately see via Corollary 4.1.10 that for all t ,

$$\mathbb{P} \left(\hat{\mu}_t(i) \geq \mu_i + \sqrt{\frac{\sigma^2 \log \frac{1}{\delta}}{N_t(i)}} \right) \vee \mathbb{P} \left(\hat{\mu}_t(i) \leq \mu_i - \sqrt{\frac{\sigma^2 \log \frac{1}{\delta}}{N_t(i)}} \right) \leq \delta. \quad (18.2.3)$$

That is, so long as we pull the arms sufficiently many times, we are unlikely to pull the wrong arm. Here then is the UCB procedure:

Input: Sub-gaussian parameter σ^2 and sequence of deviation probabilities $\delta_1, \delta_2, \dots$
Initialization: Play each arm $i = 1, \dots, K$ once
Repeat: for each iteration t , play the arm maximizing

$$\hat{\mu}_t(i) + \sqrt{\frac{\sigma^2 \log \frac{1}{\delta_t}}{N_t(i)}}.$$

Figure 18.1: The Upper Confidence Bound (UCB) Algorithm

If we define

$$\Delta_i := \mu_{i^*} - \mu_i$$

to be the gap in means between the optimal arm and any sub-optimal arm, we then obtain the following guarantee on the expected number of pulls of any sub-optimal arm i after n steps.

Proposition 18.2.2. *Assume that each of the K arms is σ^2 -sub-Gaussian and let the sequence $\delta_1 \geq \delta_2 \geq \dots$ be non-increasing and positive. Then for any T and any arm $i \neq i^*$,*

$$\mathbb{E}[N_T(i)] \leq \left\lceil \frac{4\sigma^2 \log \frac{1}{\delta_T}}{\Delta_i^2} \right\rceil + 2 \sum_{t=2}^T \delta_t.$$

Proof Without loss of generality, we assume arm 1 satisfies $\mu_1 = \max_i \mu_i$, and let arm i be any sub-optimal arm. The key insight is to carefully consider what occurs if we play arm i in the UCB procedure of Figure 18.1. In particular, if we play arm i at time t , then we certainly have

$$\hat{\mu}_t(i) + \sqrt{\frac{\sigma^2 \log \frac{1}{\delta_t}}{N_t(i)}} \geq \hat{\mu}_1(t) + \sqrt{\frac{\sigma^2 \log \frac{1}{\delta_t}}{N_1(t)}}.$$

For this to occur, at least one of the following three events must occur (we suppress the dependence on i for each of them):

$$\begin{aligned} \mathcal{E}_{1,t} &:= \left\{ \hat{\mu}_t(i) \geq \mu_i + \sqrt{\frac{\sigma^2 \log \frac{1}{\delta_t}}{N_t(i)}} \right\}, & \mathcal{E}_{2,t} &:= \left\{ \hat{\mu}_1(t) \leq \mu_1 - \sqrt{\frac{\sigma^2 \log \frac{1}{\delta_t}}{N_1(t)}} \right\}, \\ \mathcal{E}_{3,t} &:= \left\{ \Delta_i \leq 2\sqrt{\frac{\sigma^2 \log \frac{1}{\delta_t}}{N_t(i)}} \right\}. \end{aligned}$$

Indeed, suppose that none of the events $\mathcal{E}_{1,t}, \mathcal{E}_{2,t}, \mathcal{E}_{3,t}$ occur at time t . Then we have

$$\hat{\mu}_t(i) + \sqrt{\frac{\sigma^2 \log \frac{1}{\delta_t}}{N_t(i)}} < \mu_i + 2\sqrt{\frac{\sigma^2 \log \frac{1}{\delta_t}}{N_t(i)}} < \mu_i + \Delta_i = \mu_1 < \hat{\mu}_t(1) + \sqrt{\frac{\sigma^2 \log \frac{1}{\delta_t}}{N_t(1)}},$$

the inequalities following by $\mathcal{E}_{1,t}$, $\mathcal{E}_{3,t}$, and $\mathcal{E}_{2,t}$, respectively.

Now, for any $l \in \{1, \dots, n\}$, we see that

$$\begin{aligned} \mathbb{E}[N_T(i)] &= \sum_{t=1}^T \mathbb{E}[\mathbf{1}\{A_t = i\}] = \sum_{t=1}^T \mathbb{E}[\mathbf{1}\{A_t = i, N_t(i) > l\} + \mathbf{1}\{A_t = i, N_t(i) \leq l\}] \\ &\leq l + \sum_{t=l+1}^T \mathbb{P}(A_t = i, N_t(i) > l). \end{aligned}$$

Using that δ_t is non-increasing, if we set

$$l^* = \left\lceil 4 \frac{\sigma^2 \log \frac{1}{\delta_T}}{\Delta_i^2} \right\rceil,$$

then to have $N_t(i) > l^*$ it must be the case that $\mathcal{E}_{3,t}$ cannot occur at time t —that is, we would have $2\sqrt{\sigma^2 \log \frac{1}{\delta_t}}/N_t(i) > 2\sqrt{\sigma^2 \log \frac{1}{\delta_t}}/l \geq \Delta_i$. Thus we have

$$\begin{aligned} \mathbb{E}[N_T(i)] &= \sum_{t=1}^T \mathbb{E}[\mathbf{1}\{A_t = i\}] \leq l^* + \sum_{t=l^*+1}^T \mathbb{P}(A_t = i, \mathcal{E}_{3,t} \text{ fails}) \\ &\leq l^* + \sum_{t=l^*+1}^T \mathbb{P}(\mathcal{E}_{1,t} \text{ or } \mathcal{E}_{2,t}) \leq l^* + \sum_{t=l^*+1}^T 2\delta_t. \end{aligned}$$

This implies the desired result. \square

Naturally, the number of times arm i is selected in the sequential game is related to the regret of a procedure; indeed, we have

$$\text{Reg}_n = \sum_{t=1}^n (\mu_{i^*} - \mu_{A_t}) = \sum_{i=1}^K (\mu_{i^*} - \mu_i) N_i(n) = \sum_{i=1}^K \Delta_i N_i(n).$$

Using this identity, we immediately obtain two theorems on the (expected) regret of the UCB algorithm.

Theorem 18.2.3. *Let $\delta_t = \delta/t^2$ for all t . Then for any $n \in \mathbb{N}$ the UCB algorithm attains*

$$\overline{\text{Reg}}_n \leq \sum_{i \neq i^*} \frac{4\sigma^2[2\log n - \log \delta]}{\Delta_i} + \frac{\pi^2 - 2}{3} \left(\sum_{i=1}^K \Delta_i \right) \delta + \sum_{i=1}^K \Delta_i.$$

Proof First, we note that

$$\mathbb{E}[\Delta_i N_i(n)] \leq \Delta_i \left[4\sigma^2 \log \frac{1}{\delta_n} / \Delta_i^2 \right] + 2\Delta_i \sum_{t=2}^n \frac{\delta}{t^2} \leq \frac{4\sigma^2 \log \frac{1}{\delta_n}}{\Delta_i} + \Delta_i + 2\Delta_i \sum_{t=2}^n \frac{\delta}{t^2}$$

by Proposition 18.2.2. Summing over $i \neq i^*$ and noting that $\sum_{t \geq 2} t^{-2} = \pi^2/6 - 1$ gives the result. \square

Let us unpack the bound of Theorem 18.2.3 slightly. First, we make the simplifying assumption that $\delta_t = 1/t^2$ for all t , and let $\Delta = \min_{i \neq i^*} \Delta_i$. In this case, we have expected regret bounded by

$$\text{Reg}_n \leq 8 \frac{K\sigma^2 \log n}{\Delta} + \frac{\pi^2 + 1}{3} \sum_{i=1}^K \Delta_i.$$

So we see that the asymptotic regret with this choice of δ scales as $(K\sigma^2/\Delta) \log n$, roughly linear in the classes, logarithmic in n , and inversely proportional to the gap in means. As a concrete example, if we know that the rewards for each arm Y_i belong to the interval $[0, 1]$, then Hoeffding's lemma (recall Example 4.1.6) states that we may take $\sigma^2 = 1/4$. Thus the mean regret becomes at most $\sum_{i: \Delta_i > 0} \frac{2 \log n}{\Delta_i} (1 + o(1))$, where the $o(1)$ term tends to zero as $n \rightarrow \infty$.

If we knew a bit more about our problem, then by optimizing over δ and choosing $\delta = \sigma^2/\Delta$, we obtain the upper bound

$$\text{Reg}_n \leq O(1) \left[\frac{K\sigma^2}{\Delta} \log \frac{n\Delta}{\sigma^2} + K \frac{\max_i \Delta_i}{\min_i \Delta_i} \right], \quad (18.2.4)$$

that is, the expected regret scales asymptotically as $(K\sigma^2/\Delta) \log(\frac{n\Delta}{\sigma^2})$ —linearly in the number of classes, logarithmically in n , and inversely proportional to the gap between the largest and other means.

If any of the gaps $\Delta_i \rightarrow 0$ in the bound of Theorem 18.2.3, the bound becomes vacuous—it simply says that the regret is upper bounded by infinity. Intuitively, however, pulling a *slightly* sub-optimal arm should be insignificant for the regret. With that in mind, we present a slight variant of the above bounds, which has a worse scaling with n —the bound scales as \sqrt{n} rather than $\log n$ —but is independent of the gaps Δ_i .

Theorem 18.2.4. *If UCB is run with parameter $\delta_t = 1/t^2$, then*

$$\text{Reg}_n \leq \sqrt{8K\sigma^2 n \log n} + 4 \sum_{i=1}^K \Delta_i.$$

Proof Fix any $\gamma > 0$. Then we may write the regret with the standard identity

$$\text{Reg}_n = \sum_{i \neq i^*} \Delta_i N_i(n) = \sum_{i: \Delta_i \geq \gamma} \Delta_i N_i(n) + \sum_{i: \Delta_i < \gamma} \Delta_i N_i(n) \leq \sum_{i: \Delta_i \geq \gamma} \Delta_i N_i(n) + n\gamma,$$

where the final inequality uses that certainly $\sum_{i=1}^K N_i(n) \leq n$. Taking expectations with our UCB procedure and $\delta = 1$, we have by Theorem 18.2.3 that

$$\text{Reg}_n \leq \sum_{i: \Delta_i \geq \gamma} \Delta_i \frac{8\sigma^2 \log n}{\Delta_i^2} + \frac{\pi^2 + 1}{3} \sum_{i=1}^K \Delta_i + n\gamma \leq K \frac{8\sigma^2 \log n}{\gamma} + n\gamma + \frac{\pi^2 + 1}{3} \sum_{i=1}^K \Delta_i,$$

Optimizing over γ by taking $\gamma = \frac{\sqrt{8K\sigma^2 \log n}}{\sqrt{n}}$ gives the result. \square

Combining the above two theorems, we see that the UCB algorithm with parameters $\delta_t = 1/t^2$ automatically achieves the expected regret guarantee

$$\text{Reg}_n \leq C \cdot \min \left\{ \sum_{i:\Delta_i > 0} \frac{\sigma^2 \log n}{\Delta_i}, \sqrt{K\sigma^2 n \log n} \right\}. \quad (18.2.5)$$

That is, UCB enjoys some adaptive behavior. It is not, however, optimal; there are algorithms, including Audibert and Bubeck's MOSS (Minimax Optimal in the Stochastic Case) bandit procedure [11], which achieves regret

$$\text{Reg}_n \leq C \cdot \min \left\{ \sqrt{Kn}, \frac{K}{\Delta} \log \frac{n\Delta^2}{K} \right\},$$

which is essentially the bound specified by inequality (18.2.4) (which required knowledge of the Δ_i s) and an improvement by $\log n$ over the analysis of Theorem 18.2.4. It is also possible to provide a high-probability guarantee for the UCB algorithms, which follows essentially immediately from the proof techniques of Proposition 18.2.2; we leave this to Exercise 18.6.

18.3 General losses and information-based bounds

The upper confidence bound (UCB) procedure is elegant and straightforward, but in many cases we wish to move beyond the simplest K -armed bandit settings to consider more sophisticated scenarios. To that end, we now assume there is an abstract set of actions (arms) \mathcal{A} , which may or may not be finite, and we have a collection of distributions $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$ parameterized by a set Θ (abstractly, we could just let Θ be distributions, but parametric scenarios are also interesting). Often, \mathcal{A} is finite, as in the case of K -armed bandit problems, and Θ is some subset of \mathbb{R}^K (the means), but we stay in this abstract setting temporarily. We generalize things slightly to address more than just the responses $Y_t(a) \in \mathcal{Y}$, so that we have a loss function $\ell : \mathcal{A} \times \mathcal{Y} \rightarrow \mathbb{R}$ that measures the quality of an action $a \in \mathcal{A}$ for the observation $y \in \mathcal{Y}$.

Example 18.3.1 (Classical Bernoulli bandit problem): The classical bandit problem, as in the UCB case of the previous section, has actions (arms) $\mathcal{A} = \{1, \dots, K\}$, and the parameter space $\Theta = [0, 1]^K$, and P_θ is a distribution on $Y \in \{0, 1\}^K$, where Y has independent coordinates $1, \dots, K$ with $P_\theta(Y(j) = 1) = \theta_j$, that is, $Y(j) \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\theta_j)$. The goal is to find the arm $a^* \in \arg\max_j \theta_j$ with highest mean reward. Then the loss function $\ell(a, y) = -y$ satisfies $\ell(a, y) \in [-1, 0]$, and $\mathbb{E}_\theta[\ell(a, Y(a))] = -\theta_a$, so that the optimal action minimizes $\mathbb{E}_\theta[\ell(a, Y(a))]$ over actions $a \in \mathcal{A}$. \diamond

In this setting, for a given θ , we consider the expected regret

$$\text{Reg}_n(\mathcal{A}, \ell, \theta) = \mathbb{E}_\theta \left[\sum_{t=1}^n \ell(A_t, Y_t(A_t)) - \ell(A^*, Y_t(A^*)) \right], \quad (18.3.1)$$

where

$$A^* = A^*(\theta) = \underset{a \in \mathcal{A}}{\text{argmin}} \mathbb{E}_\theta[\ell(a, Y(a))]$$

minimizes the loss expected loss of taking action $a \in \mathcal{A}$ when θ is the parameter (and so is a function of θ), while $A_t \in \mathcal{A}$ is the action the player takes at time t of the process. With this abstract setting, Figure 18.2 captures the broad algorithmic framework for this section.

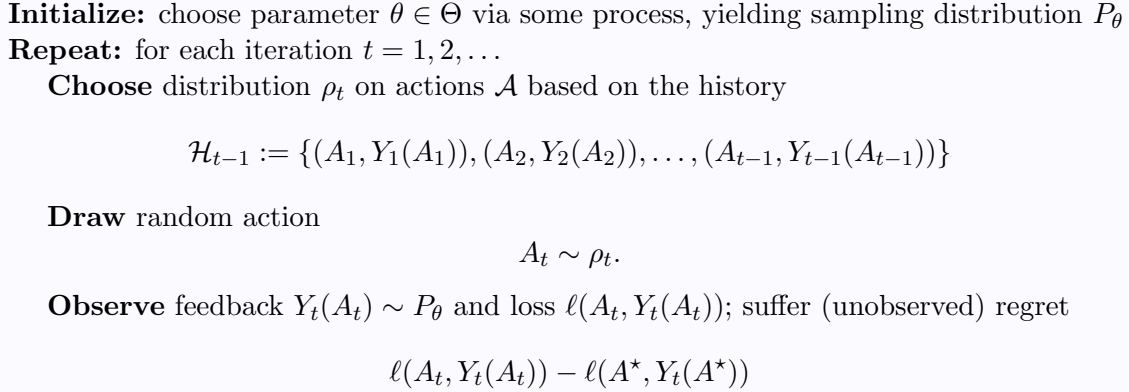


Figure 18.2: The generic exploration/exploitation algorithm

In Figure 18.2, we left the choice of θ unspecified; in worst case settings, it may be adversarial. Moving to a Bayesian setting in which $\theta \sim \pi$ for some prior distribution π on the space Θ allows us to leverage information-theoretic to obtain regret bounds, as well as new yet intuitive algorithms. Bayesian strategies—because they (can) incorporate prior knowledge—have the advantage that they suggest policies for exploration and trading between regret and information; that is, they allow us to quantify a value for information. They often yield very simple procedures, allowing simpler implementations. In this Bayesian setting where we have a prior distribution π with support Θ , we then define the Bayesian regret as

$$\text{Reg}_n(\mathcal{A}, \ell, \pi) = \mathbb{E}_\pi \left[\sum_{t=1}^n \ell(A_t, Y_t(A_t)) - \ell(A^*, Y_t(A^*)) \right], \quad (18.3.2)$$

where now the optimal action is random (as it is a function of $\theta \sim \pi$):

$$A^* \in \underset{a \in \mathcal{A}}{\text{argmin}} \mathbb{E}[\ell(a, Y(a)) \mid \theta].$$

We take the expectation (18.3.2) over both the randomness in θ according to the prior π and any randomness in the player's strategy for choosing the actions A_t at each time.

One consequence of our definition (18.3.2) worth noting is that because we take expectations, we could instead consider $\bar{\ell}(a, \theta) := \mathbb{E}[\ell(a, Y(a)) \mid \theta]$, but this quantity is typically unobservable to the algorithm. The setting (18.3.2) also encapsulates what appear to be more general settings in which we obtain observations $Y_t(a)$ by taking action $a \in \mathcal{A}$, and suffer some random loss $L_{t,a}$ that depends on the action a . But we could simply incorporate these random losses into the observation $Y_t(a)$ itself, making $\ell(a, Y_t(a))$ “pull out” the component of the observation including the loss.

18.3.1 An information-based regret bound

We can provide an information-theoretic regret bound, which bounds regret via the ratio of the expected instantaneous regret at time t and the information that taking a particular action A provides about the optimal action A^* . The rough idea is simple: if we can bound the ratio of

expected losses to the information each action A_t provides about the optimal action, then each step of the Algorithm 18.2 either gains substantial information or suffers small regret, but never both. Thus, once we have collected sufficient information, the remaining regret must be small.

To make things rigorous, we unfortunately must confront a notational choice that bedevils bounds combining information theory and statistics, in that conditional mutual information and entropy are expectations, but sometimes we wish to condition on particular realizations of variables. Thus, define the history

$$\mathcal{H}_t := \{A_1, Y_1(A_1), A_2, Y_2(A_2), \dots, A_t, Y_t(A_t)\}$$

of actions and observations through iteration t of our bandit process—here, we think of the particular realizations $(A_t, Y_t(A_t))$. Then for any action $a \in \mathcal{A}$, define the conditional mutual information

$$I_t(a) := I(A^*; Y_t(a) \mid \mathcal{H}_{t-1}) \quad (18.3.3)$$

to be the mutual information between A^* and the observation $Y_t(a)$ conditional on the realized history \mathcal{H}_{t-1} . By this we mean that we draw θ from its posterior distribution conditional on \mathcal{H}_{t-1} , then set $A^* = \operatorname{argmin}_{a \in \mathcal{A}} \ell(a, \theta)$, and draw $Y_t(a)$ conditional on θ as well. The mutual information (18.3.3) is a random variable, as it depends on the particular realization \mathcal{H}_{t-1} of the history (so we think of it analogously to $I(X; Y \mid Z = z)$ in our definitions of mutual information in Chapter 2.1.1, except the value z is random). Without any real loss of generality, we can discretize and assume \mathcal{A} is countable, allowing us to write

$$I_t(a) = H(A^* \mid \mathcal{H}_{t-1}) - H(A^* \mid Y_t(a), \mathcal{H}_{t-1}),$$

the reduction in entropy on the optimal action A^* conditional on observing $Y_t(a)$. Then for a distribution ρ on actions \mathcal{A} , define

$$I_t(\rho) := \mathbb{E}_\rho[I_t(A)] = I(A^*; Y_t(A), A \mid \mathcal{H}_{t-1}), \quad (18.3.4)$$

where now we replace the fixed action a in the definition (18.3.3) with a random $A \sim \rho$, and then conditional on $A = a$ draw $Y_t(a)$ as before. This is the information gained about A^* by taking action $A \sim \rho$, conditional on the realization \mathcal{H}_{t-1} of the history.

With the random averaged information (18.3.4), we can define the ratio between the expected regret and information gained by sampling an action $A \sim \rho$ via

$$R_t(\rho) := \frac{\mathbb{E}_\rho[\ell(A, Y_t(A)) - \ell(A^*, Y_t(A^*)) \mid \mathcal{H}_{t-1}]^2}{I_t(\rho)}. \quad (18.3.5)$$

This random ratio is, again, a random variable, as it is a function of the (realized) history \mathcal{H}_{t-1} . If we can choose a distribution ρ to make the ratio (18.3.5) small, this captures the desirable property that for $A \sim \rho$, either A should have losses close to the optimal action A^* , or A and the response $Y_t(A)$ should provide substantial information about A^* .

That gaps in squared expected losses should be related to information at all might seem *a priori* unmotivated: why the squared error? Why the information? But we have seen relationships between squared errors and information many times throughout the book; indeed, the connection between squared distances and KL-divergences underpins many of the concentration inequalities we developed in previous chapters via the Donsker-Varadhan variational representation of KL-divergence (Theorem 6.1.1). Particular examples arise in Chapter 6.2 on PAC-Bayes generalization

bounds, in our development of adaptive data analysis in Chapter 6.3 (recall Theorem 6.3.2), and in transportation inequalities for concentration (Chapter 7, Theorem 7.1.2).

By an application of the Cauchy-Schwarz inequality and using the definition of the mutual information, we can prove the next theorem, which upper bounds the expected regret by the loss/information ratio and the information gained throughout the process about the optimal arm:

Theorem 18.3.2. *For any procedure and any prior π on $\theta \in \Theta$,*

$$\text{Reg}_n(\mathcal{A}, \ell, \pi) \leq \sqrt{\mathbb{E}_\pi \left[\sum_{t=1}^n R_t(\rho_t) \right]} \sqrt{I(A^*; \{A_t, Y_t(A_t)\}_{t=1}^n)}.$$

Proof As in the discussion before the theorem, we apply the Cauchy-Schwarz inequality:

$$\begin{aligned} \mathbb{E}_\pi \left[\sum_{t=1}^n \ell(A_t, Y_t(A_t)) - \ell(A^*, Y_t(A^*)) \right] &= \mathbb{E}_\pi \left[\sum_{t=1}^n \frac{\ell(A_t, Y_t(A_t)) - \ell(A^*, Y_t(A^*))}{\sqrt{I_t(\rho_t)}} \sqrt{I_t(\rho_t)} \right] \\ &= \mathbb{E}_\pi \left[\sum_{t=1}^n \frac{\mathbb{E}_{\rho_t}[\ell(A_t, \theta) - \ell(A^*, \theta) \mid \mathcal{H}_{t-1}]}{\sqrt{I_t(\rho_t)}} \sqrt{I_t(\rho_t)} \right] \\ &= \mathbb{E}_\pi \left[\sum_{t=1}^n \sqrt{R(\rho_t)} \sqrt{I_t(\rho_t)} \right] \\ &\leq \sqrt{\mathbb{E}_\pi \left[\sum_{t=1}^n R(\rho_t) \right]} \sqrt{\mathbb{E}_\pi \left[\sum_{t=1}^n I_t(\rho_t) \right]}, \end{aligned}$$

where the second equality uses that the action sampling distribution ρ_t is a function of the history \mathcal{H}_{t-1} . Then observe that

$$\mathbb{E}_\pi[I_t(\rho_t)] = \mathbb{E}_\pi[I(A^*; A_t, Y_t(A_t) \mid \mathcal{H}_{t-1})] = I(A^*; A_t, Y_t(A_t) \mid \{A_\tau, Y_\tau(A_\tau)\}_{\tau < t})$$

by definition of the conditional mutual information. Using the chain rule for mutual information (recall Section 2.1.2),

$$\sum_{t=1}^n I(A^*; A_t, Y_t(A_t) \mid \{A_\tau, Y_\tau(A_\tau)\}_{\tau < t}) = I(A^*; \{A_t, Y_t(A_t)\}_{t=1}^n),$$

implying the theorem. □

As an immediate corollary, note that whenever the set of \mathcal{A} is finite or countable, the mutual information has the trivial upper bound

$$I(A^*; \{A_t, Y_t(A_t)\}_{t=1}^n) = H(A^*) - H(A^* \mid \{A_t, Y_t(A_t)\}_{t=1}^n) \leq H(A^*),$$

the (Shannon) entropy of A^* . We elide the dependence of A^* on the prior π over Θ , though of course different priors can yield different entropies. Whenever \mathcal{A} is finite, $H(A^*) \leq \log \text{card}(\mathcal{A})$, and we record these as a corollary.

Corollary 18.3.3. *Assume that action distributions ρ_t are chosen so that for some $R_\pi < \infty$, the average loss/information ratio satisfies*

$$\frac{1}{n} \sum_{t=1}^n \mathbb{E}_\pi[R_t(\rho_t)] \leq R_\pi$$

and that \mathcal{A} is countable. Then

$$\text{Reg}_n(\mathcal{A}, \ell, \pi) \leq \sqrt{R_\pi H(A^*)} \cdot \sqrt{n}.$$

In particular, if \mathcal{A} is finite, then

$$\text{Reg}_n(\mathcal{A}, \ell, \pi) \leq \sqrt{R_\pi \log \text{card}(\mathcal{A})} \cdot \sqrt{n}.$$

Corollary 18.3.3 thus suggests a design principal for (Bayesian) bandit algorithms: at each time step t , choose the distribution ρ_t on actions to minimize the loss/information ratio $R_t(\rho) = \mathbb{E}_\rho[\ell(A, \theta) - \ell(A^*, \theta) \mid \mathcal{H}_t]^2 / I_t(\rho)$. Though this particular choice may be impossible, it is frequently possible instead to at least guarantee that the ratio is bounded, which is evidently sufficient for a regret bound that only grows as \sqrt{n} . For the remainder of this section, we develop further consequences of these corollaries using this design principle.

18.3.2 Posterior (Thompson) sampling

One natural strategy for sampling in the Bayesian setting is to instantiate the algorithm 18.2 maintaining a posterior distribution over θ , and drawing actions that reflect this posterior. We can instantiate this via Algorithm 18.3, which describes Thompson sampling.

Input: Prior distribution π on space Θ , family of distributions $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$
Repeat: for each iteration t ,
 Choose distribution π_t to be posterior on θ given history $\mathcal{H}_{t-1} = \{A_\tau, Y_\tau(A_\tau)\}_{\tau < t}$
 Draw parameter $\theta_t \sim \pi_t$ and for $Y \sim P_{\theta_t}$, choose action

$$A_t \in \underset{a \in \mathcal{A}}{\text{argmin}} \mathbb{E}_{\theta_t}[\ell(a, Y(a))]$$

Observe feedback $Y_t(A_t)$ and loss $\ell(A_t, Y_t(A_t))$

Figure 18.3: Generic Thompson sampling algorithm

In Thompson sampling, at iteration t , we use the posterior

$$\pi_t(\theta) = \pi(\theta \mid \mathcal{H}_{t-1}),$$

the distribution on θ conditional on \mathcal{H}_{t-1} . We let

$$\rho_t^{\text{ts}} = \text{distribution of } \left\{ A_t := \underset{a \in \mathcal{A}}{\text{argmin}} \mathbb{E}[\ell(a, Y(a)) \mid \theta_t] \right\} \quad (18.3.6)$$

be the induced distribution on the actions A_t when sampling from the posterior π_t in Thompson sampling. Thompson [179] originally proposed this procedure in 1933 in the first paper on bandit problems, and it has since been the subject of substantial analysis.

We provide a few more concrete specifications of Algorithm 18.3 for Thompson (posterior), beginning with the case of Bernoulli rewards.

Example 18.3.4 (Thompson sampling for a K -armed Bernoulli bandit): Let the vector $\theta \in [0, 1]^K$ parameterize K independent $\text{Bernoulli}(\theta_a)$ distributions, where on action $a \in \mathcal{A} = \{1, \dots, K\}$, we observe $Y(a) \sim \text{Bernoulli}(\theta_a)$, that is, $\mathbb{P}(Y(a) = 1 \mid \theta) = \theta_a$.

Place a beta prior on the coordinates θ , $\theta_a \sim \text{Beta}(1, 1)$, which corresponds to the uniform distribution on $[0, 1]^d$. Let

$$N_t^1(a) = \text{card}\{\tau \leq t : A_t = a, Y_\tau(a) = 1\}$$

be the number of times arm a is pulled by time t , and similarly let $N_a^0(t) = \text{card}\{\tau \leq t : A_t = a, Y_a(\tau) = 0\}$. Then, peeking ahead to Example 19.5.2 on Beta-Bernoulli distributions, Thompson sampling with the loss $\ell(a, y) = -y$ proceeds as follows:

- (1) For each arm $a \in \mathcal{A} = \{1, \dots, K\}$, draw $\theta_t(a) \sim \text{Beta}(1 + N_a^1(t), 1 + N_a^0(t))$.
- (2) Play the action $A_t = \text{argmax}_a \theta_t(a)$.
- (3) Observe $Y_t(A_t) \in \{0, 1\}$, and increment the appropriate count.

In this case, we may implement Thompson sampling with just a few counters. \diamond

We may extend Example 18.3.4 to the case in which the losses come from any distribution with mean θ_i , so long as the distribution is supported on $[0, 1]$. In particular, we have the following example.

Example 18.3.5 (Thompson sampling with bounded random losses): Let us again consider the setting of Example 18.3.4, except that the observations $Y_t(a) \in [0, 1]$ with $\mathbb{E}[Y(a) \mid \theta] = \theta_a$. The following modification allows us to perform Thompson sampling in this case, even without knowing the distribution of $Y(a) \mid \theta$: we construct a random observation $\tilde{Y}(a) \in \{0, 1\}$ with the property that $\mathbb{P}(\tilde{Y}(a) = 1 \mid Y(a)) = Y(a)$. Then with losses $\ell(a, y) = -y$, we seek the action a maximizing the mean θ_a , and the posterior distribution over θ is still a Beta distribution. We simply redefine

$$N_a^0(t) := \text{card}\{\tau \leq t : A_t = a, \tilde{Y}_a(\tau) = 0\} \quad \text{and} \quad N_a^1(t) := \text{card}\{\tau \leq t : A_t = a, \tilde{Y}_a(\tau) = 1\}.$$

The Thompson sampling procedure is otherwise identical. \diamond

The key to our analysis of Thompson sampling is that whenever the losses themselves are sub-Gaussian, (a multiple of) the mutual information between A^* and $(Y_t(A_t), A_t)$ upper bounds the squared regret. There is some subtlety here that we note only in passing: when we say that for each $a \in \mathcal{A}$,

$$\ell(a, Y(a)) \text{ is } \sigma^2\text{-sub-Gaussian,} \tag{18.3.7}$$

we mean the following: for any distribution π on $\theta \in \Theta$, the observed loss from the process that (i) chooses $a \in \mathcal{A}$, (ii) draws $\theta \sim \pi$, and conditional on θ draws $Y \sim P_\theta$, then (iii) observes $Y(a)$ and $\ell(a, Y(a))$ is σ^2 -sub-Gaussian. Alternatively, in the context of the sequential bandit problems here, we could say that there is some σ^2 , which is a function of the history \mathcal{H}_{t-1} , such that for each $a \in \mathcal{A}$, $\ell(a, Y_t(a))$ is σ^2 -sub-Gaussian conditional on \mathcal{H}_{t-1} . A trivial sufficient condition for all of this is that the losses be bounded: if

$$\sup_y \ell(a, y) - \inf_y \ell(a, y) \leq B$$

for all $a \in \mathcal{A}$, then certainly the losses are $\frac{1}{4}B^2$ -sub-Gaussian.

The following lemma shows the information bound on the instantaneous (single step) regret. The result follows as a consequence of the Donsker-Varadhan variational representation of the KL-divergence (Theorem 6.1.1) and its connection with sub-Gaussianity.

Lemma 18.3.6. Assume that for each $a \in \mathcal{A}$, the loss $\ell(a, Y(a))$ is σ^2 -sub-Gaussian (18.3.7) and that \mathcal{A} is finite. Assume that for $\theta \sim \pi$, the optimal action $A^* = \operatorname{argmin}_{a \in \mathcal{A}} \mathbb{E}_\theta[\ell(a, Y(a))]$ has distribution ρ , and let $A \sim \rho$ independent of A^* . Then

$$\sum_{a \in \mathcal{A}} \rho(a) (\mathbb{E}[\ell(a, Y(a))] - \mathbb{E}[\ell(a, Y(a)) \mid A^* = a]) \leq \sqrt{2\sigma^2 \operatorname{card}(\mathcal{A})} \sqrt{I(A, Y(A); A^*)}.$$

Proof Let $L(a) := \ell(a, Y(a))$ be shorthand for the random realization of the loss. Then

$$\sum_{a \in \mathcal{A}} \rho(a) (\mathbb{E}[L(a)] - \mathbb{E}[L(a) \mid A^* = a]) \leq \left(\sum_a \rho(a)^2 (\mathbb{E}[L(a)] - \mathbb{E}[L(a) \mid A^* = a])^2 \right)^{1/2} \sqrt{\operatorname{card}(\mathcal{A})}.$$

Then we observe that for each $a \in \mathcal{A}$,

$$\rho(a)^2 (\mathbb{E}[L(a)] - \mathbb{E}[L(a) \mid A^* = a])^2 \leq \rho(a) \sum_{a^* \in \mathcal{A}} \rho(a^*) (\mathbb{E}[L(a)] - \mathbb{E}[L(a) \mid A^* = a^*])^2, \quad (18.3.8)$$

because we have only added nonnegative terms. Now we use that the losses are sub-Gaussian (18.3.7). If P_a and $P_{a|A^*=a^*}$ denote the marginal distribution of $L(a) = \ell(a, Y(a))$ and $L(a) = \ell(a, Y(a))$ conditional on $A^* = a^*$, respectively, then Theorem 7.1.2 implies that

$$(\mathbb{E}[L(a)] - \mathbb{E}[L(a) \mid A^* = a^*])^2 \leq 2\sigma^2 D_{\text{kl}}(P_{a|A^*=a^*} \| P_a).$$

Summing over a^* , we recognize that the marginal $P_a = \sum_{a^*} \rho(a^*) P_{a|A^*=a^*}$, and so the familiar representation (9.4.4) of the mutual information as a mixture of KL-divergences implies

$$\sum_{a^* \in \mathcal{A}} \rho(a^*) D_{\text{kl}}(P_{a|A^*=a^*} \| P_a) = I(\ell(A, Y(A)); A^* \mid A = a) \leq I(A, Y(A); A^* \mid A = a)$$

by the data processing inequality. Finally, we use the independence that A, A^* are identically distributed to observe that

$$\sum_{a \in \mathcal{A}} \rho(a) I(A, Y(A); A^* \mid A = a) = I(A, Y(A); A^*).$$

Substituting this into the bound (18.3.8), we have

$$\sum_{a \in \mathcal{A}} \rho(a)^2 (\mathbb{E}[L(a)] - \mathbb{E}[L(a) \mid A^* = a])^2 \leq 2\sigma^2 I(A, Y(A); A^*),$$

implying the lemma. \square

Lemma 18.3.6 immediately extends to the instantaneous expected loss gaps conditional on the history, because in the sampling procedure in Alg. 18.3, we have the distributional equality

$$A^* \mid \mathcal{H}_{t-1} \stackrel{\text{dist}}{=} A_t \mid \mathcal{H}_{t-1}$$

by construction of the posterior π_t on θ given the history \mathcal{H}_{t-1} . Additionally, they are certainly independent given the history, and so Lemma 18.3.6 implies that for Thompson sampling, whenever \mathcal{A} is finite,

$$(\mathbb{E}[\ell(A_t, Y_t(A_t)) \mid \mathcal{H}_{t-1}] - \mathbb{E}[\ell(A^*, Y_t(A^*)) \mid \mathcal{H}_{t-1}])^2 \leq 2\sigma^2 \operatorname{card}(\mathcal{A}) I(A_t, Y_t(A_t); A^* \mid \mathcal{H}_{t-1}).$$

Recalling the loss/information ratio (18.3.5), we see that for Thompson sampling we have the uniform bound

$$R_t(\rho_t^{\text{ts}}) \leq 2\sigma^2 \text{card}(\mathcal{A}) \quad (18.3.9)$$

uniformly for all times t . Thus, as a corollary to the main Theorem 18.3.2 bounding expected regret by the mutual information, we have the following result, which we state as a theorem to highlight the importance of Thompson sampling-style algorithms.

Theorem 18.3.7. *Let \mathcal{A} be a finite, and assume that the loss $\ell(a, Y_t(a))$ is σ^2 -sub-Gaussian (18.3.7) for each action $a \in \mathcal{A}$. Then Thompson sampling has regret*

$$\text{Reg}_n(\mathcal{A}, \ell, \pi) \leq \sqrt{2\sigma^2 \text{card}(\mathcal{A}) H(A^*)} \sqrt{n}.$$

We make a few remarks here. First, the entropy $H(A^*)$ of the optimal action is never greater than $\log \text{card}(\mathcal{A})$, and so the regret essentially scales at worst as $\sqrt{\text{card}(\mathcal{A})n}$. When the entropy of the optimal action is smaller, of course, the regret itself may be smaller. As an example application of Theorem 18.3.7, let us revisit the Bernoulli bandit problem in Example 18.3.4.

Example 18.3.8 (Bernoulli bandits, continued): Consider the Bernoulli bandit setting of Example 18.3.4. Then the mean regret for loss $\ell(a, y(a)) = -y(a)$ is

$$\text{Reg}_n = \sum_{t=1}^n \mathbb{E}_\pi[\theta_{A^*} - \theta_{A_t}],$$

the gap between the best mean reward of an arm and the played arm A_t . Because there are K arms, with $\{0, 1\}$ -valued rewards, each is $\frac{1}{4}$ -sub-Gaussian, and so Thompson sampling achieves regret

$$\text{Reg}_n \leq \sqrt{\frac{K \log K}{2}} \sqrt{n}.$$

This (up to the $\log K$ factor) minimax rate optimal, as we shall see, and removes the $\log n$ multiplicative factors (18.2.5) present in our other analyses. \diamond

JCD Comment: Add an experiment here that looks at UCB versus Thompson sampling, and see what happens.

18.3.3 Information-based exploration

The loss/information ratio (18.3.5) around which the regret bound in Theorem 18.3.2 centers suggests strategies that iteratively choose the distribution ρ to minimize the ratio. This suggests *information-directed sampling*, where we define

$$\rho_t^{\text{ids}} := \underset{\rho \in \Delta(\mathcal{A})}{\text{argmin}} \left\{ R_t(\rho) = \frac{\mathbb{E}_\rho[\ell(A, Y(A)) - \ell(A^*, Y(A^*)) \mid \mathcal{H}_{t-1}]^2}{I_t(\rho)} \right\}. \quad (18.3.10)$$

(Recall here that $\Delta(\mathcal{A})$ denotes the collection of probability distributions on \mathcal{A} , and \mathbb{E}_ρ denotes expectation over $A \sim \rho$.)

In cases where there are only finitely many actions \mathcal{A} , because Thompson sampling guarantees the bound $R_t(\rho_t) \leq 2\sigma^2 \text{card}(\mathcal{A})$ whenever the losses are sub-Gaussian (recall inequality (18.3.9)), we similarly have $R_t(\rho_t^{\text{ids}}) \leq 2\sigma^2 \text{card}(\mathcal{A})$. Thus information directed sampling achieves no worse regret:

Corollary 18.3.9. *Let ρ_t^{ids} be the information-directed sampling distribution (18.3.10), and assume the other conditions of Theorem 18.3.7. Then*

$$\text{Reg}_n(\mathcal{A}, \ell, \pi) \leq \sqrt{2\sigma^2 \text{card}(\mathcal{A}) H(A^*)} \sqrt{n}.$$

It is not always apparent how to compute the information-directed sampling distribution (18.3.10), as it involves optimization of a ratio of losses and mutual information. The connections between information and squared error, however, allow other related approaches, which can sometimes be easier to implement. For example, the following result, which is similar to Lemma 18.3.6, relates the information of actions to the variance of losses:

Lemma 18.3.10. *Assume that $\ell(a, Y(a))$ is σ^2 -sub-Gaussian (18.3.7) for each $a \in \mathcal{A}$. Then*

$$\text{Var}(\mathbb{E}[\ell(a, Y(a)) \mid A^*]) \leq 2\sigma^2 I(A^*; Y(A), A \mid A = a).$$

Proof Let $P_{a|A^*=a}$ denote the distribution of $\ell(a, Y(a))$ conditional on $A^* = a^*$, and let ρ be the distribution of A^* . Then

$$\begin{aligned} \frac{1}{2\sigma^2} I(A^*; Y(A), A \mid A = a) &\stackrel{(i)}{\geq} \frac{1}{2\sigma^2} I(A^*; \ell(a, Y(a)) \mid A = a) \\ &\stackrel{(ii)}{=} \frac{1}{2\sigma^2} \int_{a^* \in \mathcal{A}} D_{\text{kl}}(P_{a|A^*=a^*} \| P_a) d\rho(a^*) \\ &\stackrel{(iii)}{\geq} \int_{a^* \in \mathcal{A}} (\mathbb{E}[\ell(a, Y(a))] - \mathbb{E}[\ell(a, Y(a)) \mid A^* = a^*])^2 d\rho(a^*) \\ &= \text{Var}(\mathbb{E}[\ell(a, Y(a)) \mid A^*]), \end{aligned}$$

where step (i) follows by the data processing inequality, step (ii) via the representation (9.4.4) of the mutual information, and step (iii) uses the sub-Gaussianity of $\ell(a, Y(a))$ (Theorem 7.1.2), and the final equality follows because $\mathbb{E}[\mathbb{E}[\ell(a, Y(a)) \mid A^*]] = \mathbb{E}[\ell(a, Y(a))]$. \square

We consider Lemma 18.3.10 conditional on the history \mathcal{H}_{t-1} . Defining the shorthand notation

$$\mathbb{E}_t[\cdot] := \mathbb{E}[\cdot \mid \mathcal{H}_{t-1}] \quad \text{and} \quad \text{Var}_t(\cdot) := \text{Var}(\cdot \mid \mathcal{H}_{t-1})$$

for the conditional expectation and variance given the history, the lemma shows that for any distribution ρ on \mathcal{A} , the loss/information ratio satisfies $R_t(\rho) \leq 2\sigma^2 V_t(\rho)$ for the loss/variance ratio

$$V_t(\rho) := \frac{(\int_{\mathcal{A}} \mathbb{E}_t[\ell(a, Y(a)) - \ell(A^*, Y(A^*))] d\rho(a))^2}{\int_{\mathcal{A}} \text{Var}_t(\mathbb{E}_t[\ell(a, Y(a)) \mid A^*]) d\rho(a)}. \quad (18.3.11)$$

The quantity $V_t(\rho)$ is the ratio of a quadratic function in ρ over a linear function of ρ , and so is convex in ρ . (See Exercise 18.1.) Thus, at least from a computational perspective, it is natural to use *variance-directed sampling*, which chooses

$$\rho_t^{\text{vds}} := \underset{\rho \in \Delta(\mathcal{A})}{\text{argmin}} V_t(\rho).$$

Whenever we can guarantee the variance ratio $\inf_{\rho} V_t(\rho)$ is bounded for all t , then evidently we have an order \sqrt{n} regret bound.

In some cases, we can use this variance-directed sampling strategy to obtain regret bounds without having to compute information directly, instead relying on variances. For simplicity, we focus on the case when \mathcal{A} is discrete, so we can identify ρ as vectors $\rho \in \mathbb{R}^k$ or $\rho \in \mathbb{R}^N$. Then using the shorthands $v = [\text{Var}(\mathbb{E}[\ell(a, Y(a)) \mid A^*])]_{a \in \mathcal{A}}$ and $r = [\mathbb{E}[\ell(a, Y(a)) - \ell(A^*, Y(A^*))]]_{a \in \mathcal{A}}$ for the vectors of variances and instantaneous regrets, respectively, the variance ratio becomes

$$V(\rho) := \frac{(\sum_{a \in \mathcal{A}} \rho(a) \mathbb{E}[\ell(a, Y(a)) - \ell(A^*, Y(A^*))])^2}{\sum_{a \in \mathcal{A}} \rho(a) \text{Var}(\mathbb{E}[\ell(a, Y(a)) \mid A^*])} = \frac{\langle \rho, r \rangle^2}{\langle \rho, v \rangle}.$$

In the case of the K -armed bandit with sub-Gaussian arms, we can modify the argument that Thompson sampling has bounded loss/information ratio (18.3.9) to obtain an identical bound for the variance ratio:

Lemma 18.3.11. *Assume \mathcal{A} is finite. Then variance-directed sampling and Thompson sampling satisfy*

$$V_t(\rho_t^{\text{vds}}) = \inf_{\rho} V_t(\rho) \leq V_t(\rho_t^{\text{ts}}) \leq \text{card}(\mathcal{A}). \quad (18.3.12)$$

Proof By construction, the distribution ρ_t^{ts} is identical to the distribution $\rho^*(\cdot \mid \mathcal{H}_{t-1})$ of A^* , so that $\inf_{\rho \in \Delta_K} V_t(\rho) \leq V_t(\rho_t^{\text{ts}}) = V_t(\rho^*)$. We may now proceed tacitly conditioning on the history \mathcal{H}_{t-1} to avoid notational overload. For $A, A^* \stackrel{\text{iid}}{\sim} \rho^*$,

$$\begin{aligned} \mathbb{E}[\ell(A, Y(A)) - \ell(A^*, Y(A^*))] &= \sum_{a \in \mathcal{A}} \rho^*(a) \mathbb{E}[\ell(a, Y(a))] - \sum_{a \in \mathcal{A}} \rho^*(a) \mathbb{E}[\ell(a, Y(a)) \mid A^* = a] \\ &= \sum_{a \in \mathcal{A}} \rho^*(a) (\mathbb{E}[\ell(a, Y(a))] - \mathbb{E}[\ell(a, Y(a)) \mid A^* = a]), \end{aligned}$$

while the vector $v = [\text{Var}(\mathbb{E}[\ell(a, Y(a)) \mid A^*])]_{a \in \mathcal{A}}$ of variances satisfies

$$\langle \rho^*, v \rangle = \sum_{a \in \mathcal{A}} \rho^*(a) \sum_{a^* \in \mathcal{A}} \rho^*(a^*) (\mathbb{E}[\ell(a, Y(a)) \mid A^* = a^*] - \mathbb{E}[\ell(a, Y(a))])^2.$$

Jensen's inequality implies that

$$\begin{aligned} (\mathbb{E}[\ell(A, Y(A)) - \ell(A^*, Y(A^*))])^2 &\leq K \sum_{a \in \mathcal{A}} (\rho^*(a))^2 (\mathbb{E}[\ell(a, Y(a))] - \mathbb{E}[\ell(a, Y(a)) \mid A^* = a])^2 \\ &\leq K \sum_{a, a^* \in \mathcal{A}} \rho^*(a) \rho^*(a^*) (\mathbb{E}[\ell(a, Y(a))] - \mathbb{E}[\ell(a, Y(a)) \mid A^* = a^*])^2 = K \langle \rho^*, v \rangle, \end{aligned}$$

because we have simply added more positive terms. \square

As an immediate consequence of the inequality (18.3.12), we have the following corollary, which follows from the main Theorem 18.3.2.

Corollary 18.3.12. *Let ρ_t^{vds} be the variance-directed sampling distribution, and assume the other conditions of Theorem 18.3.7. Then*

$$\text{Reg}_n(\mathcal{A}, \ell, \pi) \leq \sqrt{2\sigma^2 \text{card}(\mathcal{A}) H(A^*)} \sqrt{n}.$$

18.3.4 An extended example: linear bandits

In many bandit and other exploration/exploitation problems, the rewards among arms are not independent—losses $\ell(a, Y(a))$ for a particular action $a \in \mathcal{A}$ are correlated with those for other actions $a' \in \mathcal{A}$. So-called *linear bandits* provide the cleanest model capturing this setting, making them a natural target of analysis. In this case, we identify actions \mathcal{A} with a subset of \mathbb{R}^d , and the “true” state of the world with a parameter $\theta \in \mathbb{R}^d$. Then upon taking action $a \in \mathcal{A}$, we receive feedback and loss

$$Y(a) = \langle a, \theta \rangle + \varepsilon \quad \text{and} \quad \ell(a, Y(a)) = -Y(a), \quad (18.3.13)$$

where ε is conditionally mean-zero noise and we therefore wish to maximize the inner product

$$\langle a, \theta \rangle = -\mathbb{E}[\ell(a, Y(a)) \mid \theta].$$

Example 18.3.13 (Network routing): In the source routing problem, one wishes to send a sequence of packets from a source to a given destination, and one chooses a path (sequence of links) on which to send the packets among nodes $i = 1, \dots, d$. Representing a path via a matrix $a \in \{0, 1\}^{d \times d}$, where $a_{ij} = 1$ if the packet is sent on the link from node i to j and 0 otherwise, then the (noiseless) total transit cost is $\sum_{ij} a_{ij} \theta_{ij}$, where θ_{ij} indicates the delay on edge (i, j) . The actions \mathcal{A} correspond to valid paths between nodes in the network. \diamond

Given the setting (18.3.13), we can extend Lemma 18.3.11 to apply to arbitrary action sets \mathcal{A} .

Proposition 18.3.14. *Assume the linear bandit model (18.3.13). Then variance-directed sampling and Thompson sampling satisfy*

$$V_t(\rho_t^{\text{vds}}) \leq V_t(\rho_t^{\text{ts}}) \leq d,$$

where V_t denotes the loss/variance ratio (18.3.11). If additionally $Y(a) = \langle \theta, a \rangle + \varepsilon$ is σ^2 -sub-Gaussian for each $a \in \mathcal{A}$, then information-directed sampling, variance-directed sampling, and Thompson sampling satisfy the loss/information ratio bound

$$R_t(\rho_t^{\text{ids}}) \leq 2\sigma^2 V_t(\rho_t^{\text{vds}}) \leq 2\sigma^2 V_t(\rho_t^{\text{ts}}) \leq 2\sigma^2 d.$$

Proof We tacitly ignore the history \mathcal{H}_{t-1} as it is immaterial for the proof, and let $V = V_t$ be the loss/variance ratio. Let $\rho = \rho^{\text{ts}}$, so that $A^* \sim \rho$. Then $V(\rho^{\text{vds}}) \leq V(\rho)$, and defining the matrix $M \in \mathbb{R}^{d \times d}$ by $M = \text{Cov}(\mathbb{E}[\theta \mid A^*])$, when $A \sim \rho$ the variance becomes

$$\mathbb{E}_\rho[\text{Var}(\langle \mathbb{E}[\theta \mid A^*], A \rangle \mid A)] = \langle \text{Cov}(\mathbb{E}[\theta \mid A^*]), \mathbb{E}_\rho[AA^\top] \rangle = \langle M, \mathbb{E}_\rho[AA^\top] \rangle.$$

Noting that $\mathbb{E}[A] = \mathbb{E}[A^*]$ when $A, A^* \sim \rho$, because A is independent of θ we also have

$$\begin{aligned} \mathbb{E}[\ell(A, Y(A)) - \ell(A^*, Y(A^*))] &= \mathbb{E}[\langle \theta, A^* - A \rangle] = \mathbb{E}[\langle \theta, A^* \rangle - \langle \mathbb{E}[\theta], A^* \rangle] \\ &= \mathbb{E}[\langle \mathbb{E}[\theta \mid A^*] - \mathbb{E}[\theta], A^* \rangle]. \end{aligned}$$

Then the Cauchy-Schwarz inequality implies

$$\begin{aligned} \mathbb{E}[\langle \mathbb{E}[\theta \mid A^*] - \mathbb{E}[\theta], A^* \rangle]^2 &= \mathbb{E} \left[\left\langle M^{-1/2} (\mathbb{E}[\theta \mid A^*] - \mathbb{E}[\theta]), M^{1/2} A^* \right\rangle \right]^2 \\ &\leq \mathbb{E} \left[\left\| M^{-1/2} (\mathbb{E}[\theta \mid A^*] - \mathbb{E}[\theta]) \right\|_2^2 \right] \mathbb{E} \left[\left\| M^{1/2} A^* \right\|_2^2 \right] \\ &= \langle M^{-1}, \text{Cov}(\mathbb{E}[\theta \mid A^*]) \rangle \langle M, \mathbb{E}[AA^\top] \rangle \end{aligned}$$

because A and A^* have identical distributions. Finally, we use that $\langle M^{-1}, M \rangle = \text{tr}(I_d) = d$, implying $V(\rho^{\text{vds}}) \leq V(\rho^{\text{ts}}) = V(\rho) \leq d$.

For the second claim of the proposition, we simply apply Lemma 18.3.11. \square

As a corollary, we can obtain the “standard” regret bounds for linear bandits, assuming we have a prior π over the parameters $\Theta \subset \mathbb{R}^d$. In this case, the loss

$$\ell(a, y) = -y \quad \text{satisfies} \quad \mathbb{E}[\ell(a, Y(a)) \mid \theta] = -\langle a, \theta \rangle,$$

and so minimizing the loss over $a \in \mathcal{A}$ corresponds to maximizing $\langle a, \theta \rangle$. We then obtain the following corollary.

Corollary 18.3.15. *Let \mathcal{A} be a countable or finite set of actions. Assume the linear bandit model (18.3.13) and that for each $a \in \mathcal{A}$, $Y(a) = \langle \theta, a \rangle + \varepsilon$ is σ^2 -sub-Gaussian (18.3.7). Then information-directed sampling, variance-directed sampling, and Thompson sampling each have the regret bound*

$$\text{Reg}_n(\mathcal{A}, \ell, \pi) \leq \sqrt{2\sigma^2 d H(A^*)} \cdot \sqrt{n}.$$

Let us work through one example here, which shows how Thompson-style sampling procedures can apply to linear bandit problems.

Example 18.3.16 (Thompson sampling for linear bandits with a Gaussian prior): Consider the linear bandit setting (18.3.13), and let $\Theta = \mathbb{R}^d$ be all of d -dimensional space; put a Gaussian prior $\mathbf{N}(0, \tau^2 I_d)$ on θ , where $\tau^2 < \infty$ captures the prior variance, and assume the noise $\varepsilon \stackrel{\text{iid}}{\sim} \mathbf{N}(0, \sigma^2)$. At least at the first step $t = 1$, the observed loss

$$-\ell(a, Y(a)) = Y(a) = \langle a, \theta \rangle + \varepsilon \sim \mathbf{N}\left(0, \|a\|_2^2 \tau^2 + \sigma^2\right),$$

which is sub-Gaussian so long as \mathcal{A} is bounded.

We claim that for $\mathcal{H}_t = \{Y_1(A_1), A_1, \dots, Y_t(A_t), A_t\}$, we have

$$\theta \mid \mathcal{H}_t \sim \mathbf{N}(\mathbb{E}[\theta \mid \mathcal{H}_t], t^{-1} \Sigma_t) \quad (18.3.14)$$

for

$$\Sigma_t = \left(\frac{1}{t\tau^2} I_d + \frac{1}{t\sigma^2} \sum_{i=1}^t A_i A_i^\top \right)^{-1} \quad \text{and} \quad \mathbb{E}[\theta \mid \mathcal{H}_t] = \frac{1}{\sigma^2} \Sigma_t \left(\frac{1}{t} \sum_{i=1}^t A_i Y_i(A_i) \right).$$

Assuming that the distributional equality (18.3.14) holds, then we have the necessary sub-Gaussianity (18.3.7) for the regret bounds to apply, as conditional on \mathcal{H}_t , $\ell(a, Y(a))$ is Gaussian with variance $\frac{1}{t} a^\top \Sigma_t a + \sigma^2 \leq \tau^2 \|a\|_2^2 + \sigma^2$, where we used $\Sigma_t \preceq t\tau^2$. So then for any action set \mathcal{A} and $\mathbf{N}(0, \tau^2)$ prior π , we have the following expected regret bound for information-directed, variance-directed, or Thompson sampling and linear bandits: whenever \mathcal{A} is finite,

$$\text{Reg}_n(\mathcal{A}, \ell, \pi) \leq \sqrt{2d(\tau^2 \max_{a \in \mathcal{A}} \|a\|_2^2 + \sigma^2) \log \text{card}(\mathcal{A})} \cdot \sqrt{n}.$$

Let us return to develop the posterior distribution (18.3.14) of θ . We prove the result inductively. Let $\pi(\theta \mid y_1^t, a_1^t)$ be shorthand for the density of θ given $Y_i(a_i) = y_i$ and $A_i = a_i$ for $i \leq t$. We show that

$$\pi(\theta \mid y_1^t, a_1^t) \propto \pi(\theta \mid y_1^{t-1}, a_1^{t-1}) \exp\left(-\frac{1}{2\sigma^2}(y_t - \langle \theta, a_t \rangle)^2\right). \quad (18.3.15)$$

To see this, note that at time t , we have the graphical structure in Figure 18.4, so that $Y_t(A_t)$ is conditionally independent of the history \mathcal{H}_{t-1} given A_t, θ . Thus we can write (with some abuse of notation) that

$$\begin{aligned}\pi(\theta \mid y_1^t, a_1^t) &\propto p(\theta, y_t, a_t \mid y_1^{t-1}, a_1^{t-1}) = \pi(\theta \mid y_1^{t-1}, a_1^{t-1}) p(y_t, a_t \mid \theta, y_1^{t-1}, a_1^{t-1}) \\ &= \pi(\theta \mid y_1^{t-1}, a_1^{t-1}) p(y_t \mid a_t, \theta) p(a_t \mid y_1^{t-1}, a_1^{t-1}),\end{aligned}$$

where we use p to denote an appropriate density or p.m.f. As $p(a_t \mid y_1^{t-1}, a_1^{t-1})$ does not depend on θ , we have the claim (18.3.15). For $t = 1$, we have $\pi(\theta \mid y_1^0, a_1^0) = \pi(\theta)$, and so inductively, we obtain the equality

$$\pi(\theta \mid a_1^t, y_1^t) \propto \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^t (y_i - \langle a_i, \theta \rangle)^2 - \frac{1}{2\tau^2} \|\theta\|_2^2 \right).$$

Given this equality, the conditional distribution (18.3.14) follows by algebraic manipulations (see Exercise 18.2). \diamond

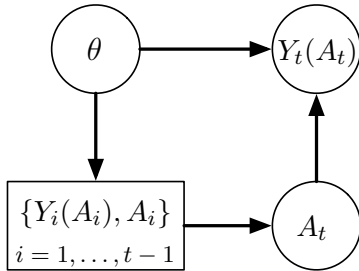


Figure 18.4. Graphical independence structure in sequential problems. Action A_t is drawn conditional on past observations and actions $\mathcal{H}_{t-1} = \{Y_i(A_i), A_i\}_{i=1}^{t-1}$, and given θ , observation $Y_t(A_t)$ is conditionally independent of \mathcal{H}_{t-1} .

The sampling distribution (18.3.14) is quite easy to use in Example 18.3.16: letting

$$\hat{\theta}_t = \operatorname{argmin} \left\{ \frac{1}{2\sigma^2} \sum_{i=1}^t (Y_i - \langle A_i, \theta \rangle)^2 + \frac{1}{2\tau^2} \|\theta\|_2^2 \right\}$$

be the minimizer of the regularized squared error, we draw $\theta_{t+1} \sim \mathbf{N}(\hat{\theta}_t, \frac{1}{t} \Sigma_t)$ for the covariance Σ_t in (18.3.14), and then choose

$$A_{t+1} \in \operatorname{argmin}_{a \in \mathcal{A}} \langle a, \theta_{t+1} \rangle.$$

For example, if $\mathcal{A} = \{-1, 1\}^d / \sqrt{d}$ consists of all $\{-1, 1\}$ -valued vectors normalized to the unit ball, then so long as $\tau^2, \sigma^2 = O(1)$, we have expected regret $\operatorname{Reg}_n \leq O(1)d\sqrt{n}$. Figure 18.5 shows example behavior with this sampling strategy, with the minor modification that we take $\mathcal{A} = \{a \in \mathbb{R}^d \mid \|a\|_2 \leq 1\}$ to be the ℓ_2 -ball. The figure shows that our analysis captures the typical behavior of the methods: the majority of the time, the (average) regret scales no worse than d/\sqrt{n} , and this appears to capture the typical behavior as well.

18.4 Online gradient descent approaches

Returning to the more basic multi-armed bandit setting of Section 18.1, it is natural to ask if we might leverage the online optimization approaches from Chapter 17 to tackle these problems.

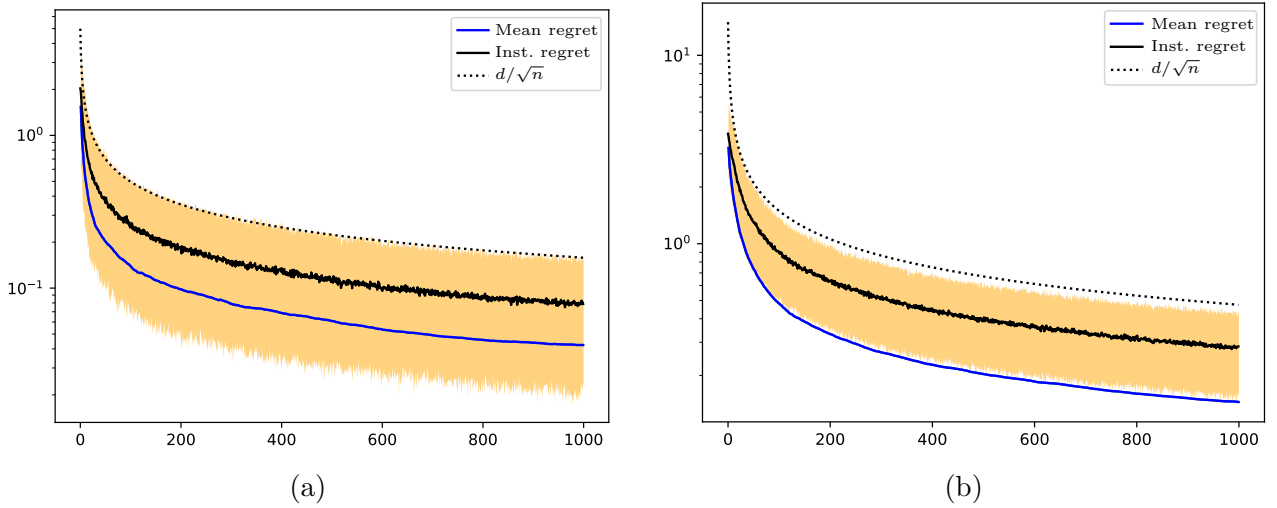


Figure 18.5. The regret behavior of Thompson sampling for the linear bandit problem (18.3.13). Each plot shows the results of 500 experiments run for 1000 iterations using the Gaussian prior of Example 18.3.16, where the action set $\mathcal{A} = \{a \in \mathbb{R}^d \mid \|a\|_2 \leq 1\}$ is the ℓ_2 -ball of radius 1. The dark black line (Inst. regret) in each plot shows the instantaneous regret $\langle \theta, A^* - A_t \rangle$ at each iteration t , averaged over the 500 experiments. The blue line (Mean regret) shows the “posterior mean” regret: taking $\bar{\theta}_t = \mathbb{E}[\theta \mid \mathcal{H}_t]$ and the “mean” action $\bar{a}_t = \bar{\theta}_t / \|\bar{\theta}_t\|_2$. The dotted black line shows d/\sqrt{n} , while the shaded orange region gives the 10–90% quantiles of the instantaneous regret across experiments.

This is certainly possible. In this scenario, we would like to formulate the bandit problem as one of minimizing a (partially) observed sequence of convex losses, where to leverage the approaches in Chapter 17, we must be able to construct unbiased gradient estimators to have any hope of achieving small regret. Let us assume as usual that we have K arms, with an (unknown) mean vector $\mu \in \mathbb{R}^K$. Then at each step t of the procedure, we play a distribution $w_t \in \Delta_K$ on the arms, and then we select an arm a at random with probability $w_{t,a}$. The *expected* loss we suffer is then $\ell_t(w_t) = \langle w_t, \mu \rangle$, though we observe only a random realization of the loss for the arm a that we play.

Because of its natural connections with estimation of probability distributions, we will use the exponentiated gradient algorithm, Example 17.2.5, to play this game. We face one main difficulty: we must estimate the gradient of the losses, $\nabla \ell_t(w_t) = \mu$, even though we only observe a random variable $Y_t(a) \in \mathbb{R}_+$, conditional on selecting action $A_t = a$ at time t , with the property that $\mathbb{E}[Y_t(a) \mid \mathcal{H}_{t-1}] = \mu_a$. Happily, we can construct such an estimate without too much additional variance.

Lemma 18.4.1. *Let $Y \in \mathbb{R}^K$ be a random variable with $\mathbb{E}[Y] = \mu$ and $w \in \Delta_K$ be a probability vector with positive entries. Choose coordinate $a \in [K]$ with probability w_a and define the random vector*

$$\tilde{Y}(j) = \begin{cases} Y(j)/w_j & \text{if } j = a \\ 0 & \text{otherwise.} \end{cases}$$

Then $\mathbb{E}[\tilde{Y} \mid Y] = Y$.

Proof The proof is immediate: for each coordinate j of \tilde{Y} , we have $\mathbb{E}[\tilde{Y}(j) \mid Y] = w_j Y(j)/w_j = Y(j)$. \square

Lemma 18.4.1 suggests the following procedure, which gives rise to (a variant of) Auer et al.'s EXP3 (Exponentiated gradient for Exploration and Exploitation) algorithm [13]. We can prove

Input: stepsize parameter η , initial vector $w_1 = [\frac{1}{K} \ \cdots \ \frac{1}{K}]^\top$
Repeat: for each iteration t ,
 Choose random action $A_t = a$ with probability $w_{t,a}$
 Receive non-negative loss $Y_t(a)$, and define

$$g_{t,j} = \begin{cases} Y_t(j)/w_j & \text{if } A_t = j \\ 0 & \text{otherwise.} \end{cases}$$

Update for each $i = 1, \dots, K$

$$w_{t+1,i} = \frac{w_{t,i} \exp(-\eta g_{t,i})}{\sum_j w_{t,j} \exp(-\eta g_{t,j})}.$$

Figure 18.6: Exponentiated gradient for bandit problems.

the following bound on the expected regret of the EXP3 Algorithm 18.6 by leveraging our refined analysis of exponentiated gradients in Proposition 17.5.1.

Proposition 18.4.2. *Assume that for each j , we have $\mathbb{E}[Y(j)^2] \leq \sigma^2$ and the observed loss $Y(j) \geq 0$. Then Alg. 18.6 attains expected regret*

$$\text{Reg}_n = \sum_{t=1}^n \mathbb{E}[\mu_{A_t} - \mu_{a^*}] \leq \frac{\log K}{\eta} + \frac{\eta}{2} \sigma^2 K n.$$

In particular, choosing $\eta = \sqrt{\log K / (K \sigma^2 n)}$ gives

$$\text{Reg}_n = \sum_{t=1}^n \mathbb{E}[\mu_{A_t} - \mu_{a^*}] \leq \frac{3}{2} \sigma \sqrt{K n \log K}.$$

Proof With Lemma 18.4.1 in place, we recall the refined regret bound of Proposition 17.5.1. For $w^* \in \Delta_K$ and any sequence of vectors g_1, g_2, \dots with $g_t \in \mathbb{R}_+^K$, exponentiated gradient descent achieves

$$\sum_{t=1}^n \langle g_t, w_t - w^* \rangle \leq \frac{\log K}{\eta} + \frac{\eta}{2} \sum_{t=1}^n \sum_{j=1}^K w_{t,j} g_{t,j}^2.$$

To transform this into a useful bound, we take expectations. Indeed, we have

$$\mathbb{E}[g_t \mid w_t] = \mathbb{E}[Y] = \mu$$

by construction, and we also have

$$\mathbb{E} \left[\sum_{j=1}^K w_{t,j} g_{t,j}^2 \mid w_t \right] = \sum_{j=1}^K w_{t,j}^2 \mathbb{E}[Y_t(j)^2 / w_{t,j}^2 \mid w_t] = \sum_{j=1}^K \mathbb{E}[Y(j)^2] = \mathbb{E}[\|Y\|_2^2].$$

This careful normalizing, allowed by Proposition 17.5.1, is essential to our analysis (and fails for more naive applications of online convex optimization bounds). In particular, we have

$$\text{Reg}_n = \sum_{t=1}^n \mathbb{E}[\langle \mu, w_t - w^* \rangle] = \sum_{t=1}^n \mathbb{E}[\langle g_t, w_t - w^* \rangle] \leq \frac{\log K}{\eta} + \frac{\eta}{2} n \mathbb{E}[\|Y\|_2^2].$$

Taking expectations gives the result. \square

Proposition 18.4.2 provides a regret bound in the multi-armed bandit problem that applies as soon as the (random) losses $Y_t(a)$ are nonnegative and have finite second moment. In this sense, it is more general than the bounds for the other methods we have developed, which we have analyzed using sub-Gaussianity. When the losses $Y(a)$ are bounded or sub-Gaussian, it is possible to achieve high-probability guarantees on the regret, though this is beyond our scope. When the random observed losses $Y_t(a)$ are bounded in $[0, 1]$, the sub-Gaussian constant $\sigma^2 = \frac{1}{4}$, yielding the mean regret bound $\frac{3}{4} \sqrt{Kn \log K}$, which is as sharp (to within constant factors) as any of our other bounds.

18.4.1 Some empirical comparisons

In spite of the similarities in the regret bounds each procedure enjoys, they frequently exhibit strikingly different performance. Practically, it appears that Thompson-sampling strategies and their variants—such as information-directed sampling—typically obtain the best empirical performance, even in non-Bayesian settings. Here, we show a few results for K -armed Bernoulli bandits, where $Y(a) \sim \text{Bernoulli}(\theta_a)$ for a vector $\theta \in [0, 1]^K$ highlighting the performance of the different methods. As a brief remark, the EXP3 algorithm (Fig. 18.6) often gets “stuck” in practice—at some point, the weights on one arm are exponentially larger than the weights on others. Thus, practically it appears important to use a variant that encourages some additional exploration by regularizing the weights, at time t defining the weights w as in Alg. 18.6, but then mixing them with a uniform distribution via

$$w_{t,a} := \frac{\exp(-\eta \sum_{\tau < t} g_{\tau,a})}{\sum_j \exp(-\eta \sum_{\tau < t} g_{\tau,j})}, \quad \epsilon_t := \min \left\{ \frac{1}{K}, \sqrt{\frac{\log K}{Kt}} \right\}, \quad \text{and} \quad \rho_t(a) = (1 - K\epsilon_t)w_{t,a} + \epsilon_t \quad (18.4.1)$$

for each $a \in \{1, \dots, K\}$. Then we draw $A_t \sim \rho_t$ in iteration t . The exploration strategy (18.4.1) encourages just enough more exploration without incurring substantial additional regret; see the references for pointers to its analysis.

In Figure 18.7, we plot results that show typical behavior of the three main algorithms we consider in this chapter for K -armed Bernoulli bandit problems: Upper Confidence Bound (UCB)-type algorithms (Section 18.2), Thompson/posterior sampling algorithms (Section 18.3.2), and ϵ -exploration variant (18.4.1) of exponentiated gradient (EXP3) algorithms (Section 18.4). Each figure plots a summary of the instantaneous regrets $\theta_{A_t} - \theta_{A^*}$ over $t = 1, \dots, 1000$ iterations, collating the results of 4000 experiments. In each individual experiment, we draw $\theta \sim \text{Uniform}([0, 1]^K)$ at the outset, then at each iteration t , conditional on the action $A_t = a \in \{1, \dots, K\}$, return $Y_t(a) \sim \text{Bernoulli}(\theta_a)$. The figure highlights a few things. First, the regrets appear to decrease polynomially in t , consistent with the analysis we have developed. Second, Thompson sampling outperforms the other algorithms; it is a reasonable default procedure when it is implementable. Finally, the bottom plot in Fig. 18.7 shows that eventually, Thompson sampling appear to find the best arm, incurring 0 regret in most samples.

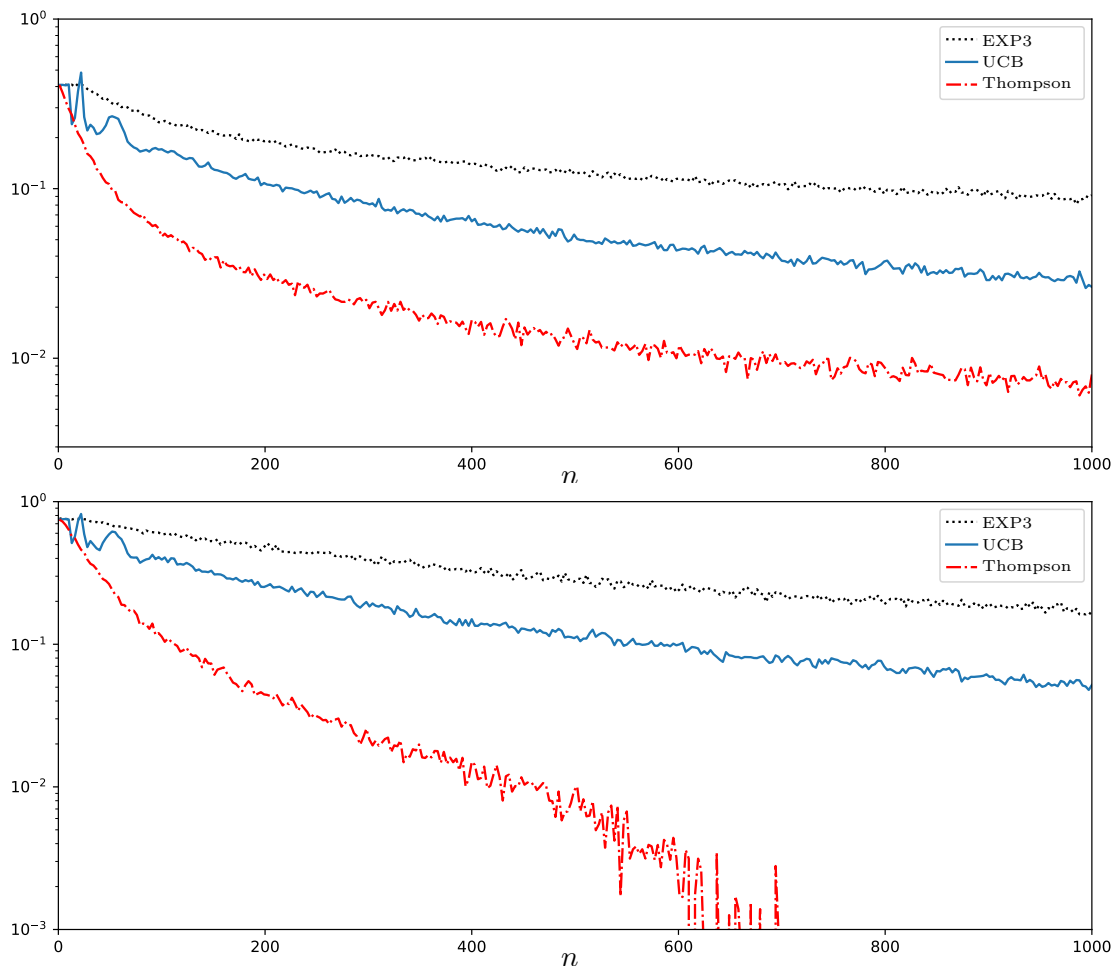


Figure 18.7. An empirical comparison of algorithms for the $K = 10$ -armed Bernoulli bandit problem. Each figure represents the results of 4000 experiments. Top: mean instantaneous regret $\theta_{A_t} - \theta_{A^*}$. Bottom: 90th percentile of instantaneous regret $\theta_{A_t} - \theta_{A^*}$ across experiments. EXP3 uses the exploration strategy (18.4.1) with $\eta = \sqrt{\log K} / \sqrt{KT}$. UCB uses confidence $\delta = 10^{-5}$. Thompson sampling uses independent $\text{Beta}(1, 1)$ priors for each coordinate of $\theta \in [0, 1]^K$.

18.5 Minimax lower bounds

Thus far, we have developed instantiations of the generic decision making framework in Figure 18.2, where we play a sequence of distributions ρ_t on actions A_t , in several scenarios, providing procedures that guarantee small regret (18.3.1) for any distribution P_θ or small Bayesian regret (18.3.2). We turn to the converse problem of lower bounds for decision-making problems. We consider two variants of the problem: a pure exploration scenario, when there is no penalty for exploration through the entire run of the procedure, and the regret based scenarios that Figure 18.2 targets. Any lower bound on the former immediately implies a lower bound on the latter: indeed, suppose that we define the action sampling distribution $\rho := \frac{1}{n} \sum_{t=1}^n \rho_t$ to be the average of the distributions

played throughout the procedure. Let $A_{n+1} \sim \rho$, and let θ parameterize P_θ . Then

$$n\mathbb{E}_\theta[\ell(A_{n+1}, Y(A_{n+1})) - \ell(A^*, Y(A^*))] = \sum_{t=1}^n \mathbb{E}_\theta[\ell(A_t, Y(A_t)) - \ell(A^*, Y(A^*))] = \text{Reg}_n(\ell, \mathcal{A}, \theta),$$

and so any lower bound on the gap

$$\mathbb{E}[\ell(A_{n+1}, Y(A_{n+1})) - \ell(A^*, Y(A^*))]$$

immediately implies one on the regret.

Let us set notation, allowing a bit more abstraction to obtain cleaner statements of the bounds. We will say P is a *model* when it maps actions $a \in \mathcal{A}$ to a distribution over losses $L \in \mathbb{R}$ and observations $Y \in \mathcal{Y}$; in the notation of Figure 18.2, we have $L = \ell(A, Y(A))$, though our abstraction allows the losses to be simply any scalar and Y to denote the observed feedback under the sampling distribution $P(a)$ upon choosing action a . Thus, we observe $(L, Y) \sim P(a)$. The *gap* of an action a for model P is

$$\text{gap}_P(a) := \sup_{a^* \in \mathcal{A}} \mathbb{E}_{P(a)}[L] - \mathbb{E}_{P(a^*)}[L], \quad (18.5.1)$$

corresponding to $\mathbb{E}[\ell(a, Y(a)) - \ell(A^*, Y(A^*))]$. We say that a distribution q is an *exploration strategy* if for each $t = 1, 2, \dots$, it defines a distribution for the t th action $A_t \sim q(\cdot \mid \mathcal{H}_{t-1})$, where $\mathcal{H}_t = \{A_i, Y_i(A_i)\}_{i \leq t}$ denotes the history as usual.

Given the gap (18.5.1) and definition of an exploration strategy, the *minimax risk* of a decision estimation problem is then

$$\mathfrak{M}_n(\mathcal{P}) := \inf_{\rho, q} \sup_{P \in \mathcal{P}} \mathbb{E}_{P, \rho, q} [\text{gap}_P(A_{n+1})], \quad (18.5.2)$$

where the infimum is taken over all exploration strategies q and distributions ρ on actions A_{n+1} , where ρ is a measurable function of $\mathcal{H}_n = (A_1^n, Y_1^n)$. In the literature on bandit problems, this quantity is the minimax regret for *best arm identification*, which cares only about identifying the optimal action $a \in \mathcal{A}$. We can define the minimax regret similarly, letting

$$\mathfrak{R}_n(\mathcal{P}) := \inf_q \sup_{P \in \mathcal{P}} \sum_{t=1}^n \mathbb{E}_{P, q} [\text{gap}_P(A_t)],$$

where we take the infimum over strategies q . The discussion above shows that

$$\mathfrak{R}_n(\mathcal{P}) \geq n \cdot \mathfrak{M}_n(\mathcal{P}),$$

so we focus essentially exclusively on $\mathfrak{M}_n(\mathcal{P})$.

18.5.1 Action separation and a modulus of continuity

To lower bound the minimax risk (18.5.2), we provide an analogue of the moduli of continuity that were central to lower bounds on estimation of individual (scalar) functionals in Chapter 13, except that now this modulus relates particularly to decision making problems. For a collection \mathcal{P} of models with associated action set \mathcal{A} , define the *action separation* of \mathcal{P}

$$\delta_{\mathcal{A}}(\mathcal{P}) := \sup_{P_1, P_2 \in \mathcal{P}} \left\{ \delta \geq 0 \mid \text{for all } a \in \mathcal{A}, \text{gap}_{P_1}(a) \geq \delta \text{ or } \text{gap}_{P_2}(a) \geq \delta \right\}.$$

This measures the separation between models in that any action $a \in \mathcal{A}$ is sub-optimal for at least one of the distributions P_1 or P_2 . For any model \bar{P} and probability distribution q on \mathcal{A} , we can define ϵ -neighborhood of \bar{P} indexed by $A \sim q$ by

$$\mathcal{P}_{q,\epsilon}(\bar{P}) := \{P \in \mathcal{P} \mid \mathbb{E}_q [D_{\text{kl}}(\bar{P}(A) \| P(A))] \leq \epsilon^2\}.$$

(To keep the notation clear, in the case that \mathcal{A} is countable, we simply mean $\mathbb{E}_q[D_{\text{kl}}(\bar{P}(A) \| P(A))] = \sum_a D_{\text{kl}}(\bar{P}(a) \| P(a)) q(a)$ to be the average KL-divergence between $\bar{P}(a)$ and $P(a)$.)

We can now define the analogue of the modulus of a parameter with respect to the Hellinger distance (13.1.1), which we shall term the *modulus of action separation* by

$$\omega_{\mathcal{A}}(\epsilon \mid \mathcal{P}, \bar{P}) := \inf_q \delta_{\mathcal{A}}(\mathcal{P}_{q,\epsilon}(\bar{P})). \quad (18.5.3)$$

Intuitively, we might expect this quantity to control lower bounds: when ϵ is small, the definition of the KL-neighborhood $\mathcal{P}_{q,\epsilon}$ means that playing actions $A \sim q$ does not allow us to determine one distribution P_1 from another P_2 in $\mathcal{P}_{q,\epsilon}$, but any action a has some non-trivial loss in at least one of P_1 and P_2 .

Continuing our analogy with Chapter 13, we have the following result, which shows that the action modulus (18.5.3) at radius $1/\sqrt{n}$ lower bounds the minimax risk:

Theorem 18.5.1. *For any model \bar{P} , not necessarily in \mathcal{P} , the minimax risk satisfies*

$$\mathfrak{M}_n(\mathcal{P}) \geq \frac{1}{4} \omega_{\mathcal{A}}\left(\frac{1}{\sqrt{8n}} \mid \mathcal{P}, \bar{P}\right).$$

We defer the proof of Theorem 18.5.1 to Section 18.5.3, instead focusing on some applications of the result here.

The key in Theorem 18.5.1 is that a geometric-like quantity—the modulus (18.5.3)—guarantees lower bounds on the minimax risk for decision estimation. As a consequence, if we can provide lower bounds on the modulus, we immediately lower bounds on the minimax risk. The next two examples exhibit recipes for this strategy:

Example 18.5.2 (Multi-armed bandits): By Theorem 18.5.1, to prove a lower bound we need only show that a family \mathcal{P} of models, restricted to a particular neighborhood of some model \bar{P} , has large enough action separation. Consider the k -armed Gaussian bandit, where we define the model family to be collections of normal distributions with $P(a) = \mathcal{N}(\mu_a, \sigma^2)$ for arbitrary mean vectors $\mu \in \mathbb{R}^k$. Then for the “null” model \bar{P} with $\bar{P}(a) = \mathcal{N}(0, \sigma^2)$, for any p.m.f. q and a model with associated mean reward $\mu \in \mathbb{R}^k$, we have

$$\sum_{i=1}^k q_i D_{\text{kl}}(\bar{P}(i) \| P(i)) = \frac{1}{2\sigma^2} \sum_{i=1}^k q_i \mu_i^2.$$

So when $\mu = \delta e_a$ for a standard basis vector e_a , we obtain $\mathbb{E}_q[D_{\text{kl}}(\bar{P}(A) \| P(A))] = \frac{q_a \delta^2}{2\sigma^2}$.

We can now compute a lower bound on the action separation. For any q , there is necessarily at least one index i with $q_i \leq \frac{1}{k}$, and at least one distinct index j such that $q_j \leq \frac{1}{k-1}$. (Otherwise, we would have $\sum_i q_i > 1$.) Now, let P_1 and P_2 correspond to mean vectors $\mu = \delta e_i$ and δe_j for these two indices. Then it is immediate that at least one of $\text{gap}_{P_1}(a) \geq \delta$ or $\text{gap}_{P_2}(a) \geq \delta$, while

$$\mathbb{E}_q[D_{\text{kl}}(\bar{P}(A) \| P_1(A))] \leq \frac{\delta^2}{2k\sigma^2} \quad \text{and} \quad \mathbb{E}_q[D_{\text{kl}}(\bar{P}(A) \| P_2(A))] \leq \frac{\delta^2}{2(k-1)\sigma^2}.$$

So for any distribution q on A and $\epsilon > 0$, we obtain action separation

$$\delta_{\mathcal{A}}(\mathcal{P}_{q,\epsilon}(\bar{P})) \geq \sup \left\{ \delta \geq 0 \mid \frac{\delta^2}{2(k-1)\sigma^2} \leq \epsilon^2 \right\} = \sigma\epsilon\sqrt{2(k-1)}.$$

In particular, $\omega_{\mathcal{A}}(\epsilon \mid \mathcal{P}, \bar{P}) \geq \sigma\epsilon\sqrt{2(k-1)}$, and so the minimax risk for best-arm identification in a k -armed Gaussian bandit is at least

$$\mathfrak{M}_n(\mathcal{P}) \geq \frac{\sigma\sqrt{k-1}}{8\sqrt{n}}$$

by Theorem 18.5.1. As in our discussion of UCB algorithms at the end of Section 18.2, this is sharp. \diamond

Example 18.5.3 (Linear bandits): Consider linear bandits, where we index models P by $\theta \in \mathbb{B}_2^d$, take actions $x \in \mathcal{A} := \mathbb{B}_2^d$, and observe reward

$$y(x) = \langle x, \theta \rangle + \varepsilon, \quad \varepsilon \stackrel{\text{iid}}{\sim} \mathbf{N}(0, \sigma^2).$$

This is a standard linear regression problem where we can design x . For a given model P_θ , we have $\sup_{\|x\|_2 \leq 1} \mathbb{E}_{P_\theta}[y(x)] = \|\theta\|_2$, and so the gap

$$\text{gap}_{P_\theta}(x) = \|\theta\|_2 - \langle x, \theta \rangle.$$

For any distribution q on actions x , if we take $\bar{P}(x) = \mathbf{N}(0, \sigma^2)$ independent of x , we obtain

$$\mathbb{E}_q[D_{\text{kl}}(\bar{P}(X) \| P_\theta(X))] = \frac{1}{2\sigma^2} \mathbb{E}_q[\langle \theta, X \rangle^2].$$

Let $C_q := \mathbb{E}_q[XX^\top]$, where $\text{tr}(C_q) = \mathbb{E}[\|X\|_2^2] \leq 1$ and therefore there must exist distinct eigenvectors u, v of C_q for which $u^\top C_q u \leq \frac{1}{d}$ and $v^\top C_q v \leq \frac{1}{d-1}$.

Fixing $\delta \geq 0$ to be chosen, take P_1 to be the linear bandit model indexed by $\theta = \delta u$ and P_2 that indexed by $\theta = \delta v$, so that

$$\mathbb{E}_q[D_{\text{kl}}(\bar{P}(X) \| P_1(X))] \leq \frac{\delta^2}{2d\sigma^2} \quad \text{and} \quad \mathbb{E}_q[D_{\text{kl}}(\bar{P}(X) \| P_2(X))] \leq \frac{\delta^2}{2(d-1)\sigma^2}.$$

Moreover, for any $\|x\|_2 \leq 1$, we have

$$\text{gap}_{P_1}(x) = \delta(1 - \langle x, u \rangle) \quad \text{and} \quad \text{gap}_{P_2}(x) = \delta(1 - \langle x, v \rangle).$$

Because $x = \frac{v+u}{\|v+u\|_2}$ minimizes $\max\{-\langle x, v \rangle, -\langle x, u \rangle\}$ and $\|u+v\|_2 = \sqrt{2}$ by the orthogonality of u and v , we obtain

$$\max\{\text{gap}_{P_1}(x), \text{gap}_{P_2}(x)\} \geq \delta \left(1 - \frac{1}{\sqrt{2}}\right) = \frac{2 - \sqrt{2}}{2} \delta.$$

In particular, we have action separation

$$\delta_{\mathcal{A}}(\mathcal{P}_{q,\epsilon}(\bar{P})) \geq \frac{2 - \sqrt{2}}{2} \sup \left\{ \delta \geq 0 \mid \frac{\delta^2}{2(d-1)\sigma^2} \leq \epsilon^2 \right\} = (\sqrt{2} - 1) \cdot \sigma\epsilon\sqrt{d-1}.$$

The minimax rate for dimension $d \geq 1$ therefore satisfies

$$\mathfrak{M}_n(\mathcal{P}) \gtrsim \frac{\sigma\sqrt{d}}{\sqrt{n}},$$

which follows by taking $\epsilon = O(1)/\sqrt{n}$. \diamond

This minimax rate is not quite sharp: the correct rate is of order $d\sigma/\sqrt{n}$, which motivates a different lower bounding technique we explore in the next section; the application to standard multi-armed bandits, however, highlights the applications of the approach.

18.5.2 Assouad's method for lower bounds

Assouad's method, which we initially develop in Chapter 9.5, extends elegantly to adaptive problems. The key insight is that when we reduce to collections of binary tests, as in the lower bound (9.5.4), we can apply chain rules for KL-divergence even with adaptive information gathering. Here, we apply this idea to two linear bandit problems, which helps both to highlight the ideas behind its application and in providing sharp minimax and Bayesian lower bounds for both pure exploration problems and regret.

Recall the key steps in applying Assouad's method: we embed the problem in the hypercube $\mathcal{V} = \{\pm 1\}^d$, demonstrate a loss with a Hamming separation, and then show that distributions $P_v, P_{v'}$ for which v, v' differ in only a single coordinate are close in KL-divergence. Following this recipe, we will work in the linear bandit with Gaussian noise setting, where for a parameter $\theta \in \mathbb{R}^d$ and actions $x \in \mathbb{R}^d$, we observe rewards

$$y(x) = \langle \theta, x \rangle + \varepsilon \quad \text{and} \quad \text{gap}_{P_\theta}(a) = \sup_{x \in \mathcal{A}} \langle \theta, x \rangle - \langle \theta, a \rangle.$$

The interactions between the action set \mathcal{A} and the parameter space Θ are sophisticated, so we focus on two simple but natural cases: we always take $\Theta = \mathbb{B}_2^d$, the ℓ_2 -ball, and consider either the actions $\mathcal{A} = [-1, 1]^d$ or $\mathcal{A} = \mathbb{B}_2^d$. We begin by stating the main theorem of this section and a few corollaries. In the theorem, we elaborate the notation (18.5.2) for the minimax risk in a decision estimation problem by also including the action set \mathcal{A} , so $\mathfrak{M}_n(\mathcal{P}, \mathcal{A}) = \inf_{\rho, q} \sup_{P \in \mathcal{P}} \mathbb{E}_{P, \rho, q}[\text{gap}_P(A_{n+1})]$ and ρ and q draw actions in \mathcal{A} .

Theorem 18.5.4. *Let \mathcal{P} be the collection of linear bandit models $y(x) = \langle \theta, x \rangle + \varepsilon$, $\|\theta\|_2 \leq 1$, with $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ -noise. Then the following minimax lower bounds hold.*

(i) *For any action set satisfying $\{-1, 1\}^d \subset \mathcal{A} \subset [-1, 1]^d$,*

$$\mathfrak{M}_n(\mathcal{P}, \mathcal{A}) \geq \frac{1}{4} \cdot \left(\frac{d\sigma}{\sqrt{n}} \wedge \sqrt{d} \right).$$

(ii) *For any action set satisfying $\{\pm 1/\sqrt{d}\}^d \subset \mathcal{A} \subset \mathbb{B}_2^d$,*

$$\mathfrak{M}_n(\mathcal{P}, \mathcal{A}) \geq \frac{1}{8} \cdot \left(\frac{d\sigma}{\sqrt{n}} \wedge 1 \right).$$

Before proving the theorem, we can develop a few simple corollaries to the result and proof technique by considering Bayesian regret (18.3.2).

Corollary 18.5.5. *Let the conditions of Theorem 18.5.4 hold and assume that $n \geq d$. Let \mathcal{A} either satisfy condition (i) or (ii) of the theorem and let $\ell_{\text{id}}(a, y) = y$ be the identity loss. Then there is a $\delta \in [0, 1/\sqrt{d}]$ such that for π_{Uniform} uniform over $\theta \in \{-\delta, \delta\}^d$,*

$$\text{Reg}_n(\ell_{\text{id}}, \mathcal{A}, \pi_{\text{Uniform}}) \geq \frac{1}{8} d\sigma\sqrt{n}.$$

Of course, the worst-case regret also satisfies the lower bound in the corollary, so that

$$\sup_{\theta \in \Theta} \text{Reg}_n(\ell_{\text{id}}, \mathcal{A}, \theta) \geq \frac{1}{8} d \sigma \sqrt{n}$$

under the same conditions. The corollary shows that our analyses of information-directed and Thompson sampling were sharp, at least to within numerical constants; Corollary 18.3.15 shows that if $\mathcal{A} = \{\pm 1/\sqrt{d}\}^d$ or $\mathcal{A} = \{-1, 1\}^d$, then because $H(A^*) \leq d \log 2$, we have

$$\text{Reg}_n(\mathcal{A}, \ell_{\text{id}}, \pi) \leq \sqrt{2 \log 2} \cdot d \sigma \sqrt{n}$$

for any prior π .

Theorem 18.5.4 and Corollary 18.5.5 also highlight how the action modulus (18.5.3), while successful for K -armed bandit problems, fails to capture the correct lower bounds for linear bandit problems. Exercise 18.3 asks you to show how to use stochastic subgradient methods to achieve the lower bounds in Theorem 18.5.4 for the pure exploration case, highlighting a fairly straightforward approach. Exercise 18.4 shows that a minor variant of the problem leads to quite different convergence behavior, highlighting some of the subtleties in decision estimation problems, especially with regards to the interaction between action sets \mathcal{A} and the underlying distributions.

Proof of Theorem 18.5.4. We now return to prove Theorem 18.5.4 by following the Assouad's method recipe. Beginning with the embedding, for $v \in \mathcal{V} = \{-1, 1\}^d$, let $\theta_v = \frac{\delta}{\sqrt{d}} v$, where $\delta \leq 1$ is to be chosen in each lower bound construction. Then next lemma then demonstrates a Hamming separation (9.5.1) for each action set.

Lemma 18.5.6. *Let $\theta \in \{\theta_v\}_{v \in \mathcal{V}} = \{\pm \delta/\sqrt{d}\}^d$ as above, and let $a^*(\theta) = \arg\max_{a \in \mathcal{A}} \langle a, \theta \rangle$. Then*

(i) *For any action set satisfying $\{-1, 1\}^d \subset \mathcal{A} \subset [-1, 1]^d$,*

$$\langle \theta, a^*(\theta) - a \rangle \geq \frac{\delta}{\sqrt{d}} \sum_{j=1}^d \mathbf{1} \{ \text{sign}(a_j) \neq \text{sign}(\theta_j) \}.$$

(ii) *For any action set satisfying $\{\pm 1/\sqrt{d}\}^d \subset \mathcal{A} \subset \mathbb{B}_2^d$,*

$$\langle \theta, a^*(\theta) - a \rangle \geq \frac{\delta}{2d} \sum_{j=1}^d \mathbf{1} \{ \text{sign}(a_j) \neq \text{sign}(\theta_j) \}.$$

Proof For the first claim of the lemma, we have $a^*(\theta) = \text{sign}(\theta)$ (defined elementwise), and so

$$\langle \theta, a^*(\theta) - a \rangle = \sum_{j=1}^d (|\theta_j| - a_j \theta_j) \geq \sum_{j=1}^d |\theta_j| \mathbf{1} \{ \text{sign}(a_j) \neq \text{sign}(\theta_j) \},$$

where we used that $a_j \in [-1, 1]$. Noting that $|\theta_j| = \delta/\sqrt{d}$ gives the first result.

For the second result, note that whenever $\|a\|_2 \leq 1$ and $v \in \{-1, 1\}^d$, we have

$$\sqrt{d} \sum_{j=1}^d \left(\frac{1}{\sqrt{d}} - a_j v_j \right)^2 = \sqrt{d} (1 - 2 \langle a, v \rangle / \sqrt{d} + \|a\|_2^2) \leq 2\sqrt{d} \left(1 - \langle a, v \rangle / \sqrt{d} \right).$$

Scaling terms appropriately, for $v = \sqrt{d}\theta/\delta$ we obtain

$$\begin{aligned} \frac{2}{\delta} \langle \theta, a^*(\theta) - a \rangle &= \frac{2}{\sqrt{d}} \langle v, a^*(\theta) - a \rangle = 2 \left(1 - \langle v, a \rangle / \sqrt{d} \right) \\ &\geq \sum_{j=1}^d \left(\frac{1}{\sqrt{d}} - a_j v_j \right)^2 \geq \sum_{j=1}^d \frac{1}{d} \mathbf{1} \{ \text{sign}(a_j) \neq v_j \}. \end{aligned}$$

Multiplying by $\frac{\delta}{2}$ gives the lemma. \square

As the last step, we must control the variation distance terms in the lower bound (9.5.4). Let $P_{v,+j}$ and $P_{v,-j}$ denote the distribution of the observations Y_1, \dots, Y_n under the model $\theta = \delta v / \sqrt{d}$ where v_j is set to (respectively) 1 or -1 . Then we have the following lemma:

Lemma 18.5.7. *Let $\mathcal{A} \subset \mathbb{R}^d$ be any action set satisfying $\|a\|_2 \leq D < \infty$ for all $a \in \mathcal{A}$. Then for the Gaussian linear bandit model,*

$$\frac{1}{d2^d} \sum_{j=1}^d \sum_{v \in \mathcal{V}} \|P_{v,+j} - P_{v,-j}\|_{\text{TV}}^2 \leq \frac{n\delta^2}{4\sigma^2 d^2} D^2.$$

Proof Fix $v \in \{-1, 1\}^d$ temporarily, and let $P = P_{v,+j}$ and $Q = P_{v,-j}$. Use the shorthand that P_i is the distribution of the i th observation Y_i conditional on the history \mathcal{H}_{i-1} and similarly for Q_i . Let $\theta_+ \in \mathbb{R}^d$ and $\theta_- \in \mathbb{R}^d$ be the parameters associated with $P_{v,+j}$ and $P_{v,-j}$, respectively. Then applying the chain rule for the KL-divergence (Lemma 2.1.9), we obtain

$$\begin{aligned} D_{\text{kl}}(P_{v,+j} \| P_{v,-j}) &= \sum_{i=1}^n \mathbb{E}_P [D_{\text{kl}}(P_i \| Q_i)] = \sum_{i=1}^n \mathbb{E}_P [D_{\text{kl}}(\mathbf{N}(\langle A_i, \theta_+ \rangle, \sigma^2) \| \mathbf{N}(\langle A_i, \theta_- \rangle, \sigma^2) \mid A_i)] \\ &= \sum_{i=1}^n \mathbb{E}_P \left[\frac{\langle A_i, \theta_+ - \theta_- \rangle^2}{2\sigma^2} \right] = \frac{\delta^2}{2d\sigma^2} \sum_{i=1}^n \mathbb{E}[A_{ij}^2]. \end{aligned}$$

Now averaging over all $v \in \{\pm 1\}^d$ and coordinates $j = 1, \dots, d$, we obtain

$$\frac{1}{2^d d} \sum_{j=1}^d \sum_{v \in \mathcal{V}} D_{\text{kl}}(P_{v,+j} \| P_{v,-j}) \leq \frac{\delta^2}{2\sigma^2 d^2} \sum_{i=1}^n \mathbb{E}[\|A_i\|_2^2],$$

and applying Pinsker's inequality gives the result. \square

Combining Lemmas 18.5.6 and 18.5.7, we can complete the proof of the theorem by considering the cases (i) and (ii) distinguishing the action sets \mathcal{A} . We first consider case (i), the hypercube-like action set, where $\{-1, 1\}^d \subset \mathcal{A} \subset [-1, 1]^d$. Then part (i) of Lemma 18.5.6, As-souad's method (9.5.4), and Lemma 18.5.7 (take $D = \sqrt{d}$) show that

$$\mathfrak{M}_n(\mathcal{P}) \geq \frac{\sqrt{d}\delta}{2} \left[1 - \left(\frac{n\delta^2}{4\sigma^2 d} \right)^{1/2} \right].$$

Setting $\delta = \sqrt{d\sigma^2/n} \wedge 1$ implies the first claim of the theorem.

For case (ii), the ℓ_2 -type action set satisfying $\{\pm 1/\sqrt{d}\}^d \subset \mathcal{A} \subset \mathbb{B}_2^d$, the same argument (except using part (ii) of Lemma 18.5.6 and setting the action radius $D = 1$) implies

$$\mathfrak{M}_n(\mathcal{P}) \geq \frac{\delta}{4} \left[1 - \left(\frac{n\delta^2}{4\sigma^2 d^2} \right)^{1/2} \right].$$

Set $\delta = \frac{d\sigma}{\sqrt{n}} \wedge 1$.

18.5.3 Proof of Theorem 18.5.1

We begin with a minor remark to avoid some trivialities. If $\omega_{\mathcal{A}}(\epsilon \mid \mathcal{P}, \bar{P}) > 0$, then *a fortiori* the set $\mathcal{P}_{q,\epsilon}(\bar{P})$ is non-empty for all distributions q on actions, because $\delta_{\mathcal{A}}(\emptyset) = 0$. We thus proceed as if $\mathcal{P}_{q,\epsilon}(\bar{P})$ is non-empty, as otherwise the lower bound of 0 is vacuous. We shall also assume w.l.o.g. that \mathcal{A} is discrete, though this is only for notational convenience.

We begin by letting $\delta > 0$ to be chosen (we shall chose δ eventually to scale as $\omega_{\mathcal{A}}(\epsilon \mid \mathcal{P}, \bar{P})$ for $\epsilon \lesssim 1/\sqrt{n}$). For simplicity in notation, we assume the observations Y include the observed losses. For any exploration strategy q and probability distribution ρ , the latter of which is defined conditional on the history $\mathcal{H}_n := (Y_1^n, A_1^n)$, we can define the distribution $\mathbb{P}_{P,\rho,q}$ on actions by $\mathbb{P}_{P,\mathcal{A},q}(A_{n+1} = a) = \mathbb{E}_{P,q}[\rho(a \mid \mathcal{H}_n)]$, where the expectation is taken jointly over Y_1^n drawn from P under the exploration strategy q . We also define

$$q_P(a) := \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{P,q}[q(a \mid \mathcal{H}_{i-1})] \quad (18.5.4)$$

to be the average exploration distribution.

By Markov's inequality, for any model P and distributions ρ, q , we have the trivial inequality

$$\mathbb{E}_{P,\rho,q}[\text{gap}_P(A)] \geq \delta \cdot \mathbb{P}_{P,\rho,q}(\text{gap}_P(A) \geq \delta).$$

We modify the probabilistic model to an arbitrary \bar{P} to allow an easier testing lower bound, observing that

$$\mathbb{E}_{P,\rho,q}[\text{gap}_P(A)] \geq \delta \cdot \left(\mathbb{P}_{\bar{P},\rho,q}(\text{gap}_P(A) \geq \delta) - \left\| \mathbb{P}_{\bar{P},\rho,q} - \mathbb{P}_{P,\rho,q} \right\|_{\text{TV}} \right).$$

As is typical in our lower bound proofs, we now show that the first probability is constant, and use that P and \bar{P} are “close” and a tensorization argument to bound the variation distance. Let $\mathbb{P}_{P,q}$ be the joint distribution of the observations (Y_1^n, A_1^n) under the model P and q , so that because we have the probabilistic structure $(A_1^n, Y_1^n) \rightarrow A_{n+1}$, the data processing inequality implies

$$\left\| \mathbb{P}_{\bar{P},\rho,q} - \mathbb{P}_{P,\rho,q} \right\|_{\text{TV}}^2 \leq \left\| \mathbb{P}_{\bar{P},q} - \mathbb{P}_{P,q} \right\|_{\text{TV}}^2 \leq \frac{1}{2} D_{\text{kl}} \left(\mathbb{P}_{\bar{P},q} \parallel \mathbb{P}_{P,q} \right)$$

by Pinsker's inequality.

Recognizing that at step i under exploration distribution q and model P , Y_i has distribution $Y_i \mid A_i = a \sim P(a)$ and that, by the Markovian action sampling structure that, conditional on the history, $A_i \mid \mathcal{H}_{i-1} \sim q(\cdot \mid \mathcal{H}_{i-1})$ under both \bar{P} and P , Lemma 2.1.9 then implies that

$$\begin{aligned} D_{\text{kl}} \left(\mathbb{P}_{\bar{P},q} \parallel \mathbb{P}_{P,q} \right) &= \sum_{i=1}^n \mathbb{E}_{\bar{P},q} \left[D_{\text{kl}} \left(\mathbb{P}_{\bar{P},q}(A_i \in \cdot, Y_i \in \cdot \mid \mathcal{H}_{i-1}) \parallel \mathbb{P}_{P,q}(A_i \in \cdot, Y_i \in \cdot \mid \mathcal{H}_{i-1}) \right) \right] \\ &= \sum_{i=1}^n \mathbb{E}_{\bar{P},q} \left[D_{\text{kl}} \left(\bar{P}(A_i) \parallel P(A_i) \right) \right] = n \mathbb{E}_{q_{\bar{P}}} \left[D_{\text{kl}} \left(\bar{P}(A) \parallel P(A) \right) \right], \end{aligned}$$

where the latter equality follows by definition (18.5.4) of q_P .

Combining the steps, we have now shown that for any pair ρ, q and any models P, \bar{P} , we have

$$\mathbb{E}_{P, \rho, q}[\text{gap}_P(A_{n+1})] \geq \delta \left(\mathbb{P}_{\bar{P}, \rho, q}(\text{gap}_P(A_{n+1}) \geq \delta) - \sqrt{\frac{n}{2} \mathbb{E}_{q_{\bar{P}}} [D_{\text{kl}}(\bar{P}(A) \| P(A))]} \right).$$

Now we use the definition of the action separation: for any distribution q on actions A , we have $\inf_{q^*} \delta_{\mathcal{A}}(\mathcal{P}_{q^*, \epsilon}(\bar{P})) \leq \delta_{\mathcal{A}}(\mathcal{P}_{q, \epsilon}(\bar{P}))$, and in particular, this holds for $q = q_{\bar{P}}$. So taking $\delta < \delta_{\mathcal{A}}(\mathcal{P}_{q_{\bar{P}}, \epsilon}(\bar{P}))$, there must exist P_1 and P_2 so that for all $a \in \mathcal{A}$ we have at least one of $\text{gap}_{P_1}(a) \geq \delta$ or $\text{gap}_{P_2}(a) \geq \delta$ and $\mathbb{E}_{q_{\bar{P}}} [D_{\text{kl}}(\bar{P}(a) \| P_i(a))] \leq \epsilon^2$ for each $i = 1, 2$. In particular,

$$\begin{aligned} & 2 \max_{i \in \{1, 2\}} \mathbb{E}_{P_i, \rho, q}[\text{gap}_{P_i}(A_{n+1})] \\ & \geq \delta \left(\mathbb{P}_{\bar{P}, \rho, q}(\text{gap}_{P_1}(A_{n+1}) \geq \delta) + \mathbb{P}_{\bar{P}, \rho, q}(\text{gap}_{P_2}(A_{n+1}) \geq \delta) - 2\sqrt{\frac{n}{2}\epsilon^2} \right) \\ & \geq \delta \left(1 - \sqrt{2n\epsilon^2} \right), \end{aligned}$$

because at least one of the events $\text{gap}_{P_i}(A_{n+1}) \geq \delta$ must occur.

Finally, we substitute appropriate values. We required only that $\delta < \delta_{\mathcal{A}}(\mathcal{P}_{q_{\bar{P}}, \epsilon}(\bar{P}))$, so we may take it up to $\omega_{\mathcal{A}}(\epsilon \mid \mathcal{P}, \bar{P})$. So long as the separation $\epsilon^2 \leq \frac{1}{8n}$, we have $2n\epsilon^2 \leq \frac{1}{4}$, and therefore

$$\sup_{P \in \mathcal{P}} \mathbb{E}_{P, \rho, q}[\text{gap}_P(A_{n+1})] \geq \frac{\delta}{2} \left(1 - \sqrt{2n\epsilon^2} \right) \geq \frac{\delta}{4} = \frac{1}{4} \omega_{\mathcal{A}}(1/\sqrt{8n} \mid \mathcal{P}, \bar{P})$$

as desired.

18.6 Technical proofs

18.6.1 Proof of Lemma 18.2.1

Without loss of generality, we can assume that any algorithm choosing arms to pull is, at time t , a function of the arms A_1, \dots, A_t , responses $Y_1(A_1), \dots, Y_t(A_t)$, and an auxiliary random variable $U \sim \text{Uniform}[0, 1]$. Let $\mathcal{H}_t = \{A_1, \dots, A_t, Y_t(A_1), \dots, Y_t(A_t), U\}$ be the history up to time t (more rigorously, the σ -field generated by the actions and responses to time t , along with the randomness U). Note that without loss of generality, we have $A_t \in \mathcal{H}_{t-1}$, as A_t is a function of this history and the randomness U in the algorithm.

To see this, we use the standard fact that the characteristic function of a random variable completely characterizes the random variable. Let $\varphi_i(\lambda) = \mathbb{E}[e^{\lambda Y(i)}] = \mathbb{E}[e^{\iota \lambda Y'(i)}]$, where $\iota = \sqrt{-1}$ is the imaginary unit, denote the characteristic function of $Y(i)$, which is identical to that of $Y'(i)$ by construction. Then writing the joint characteristic function of $N_t(i) \hat{\mu}_t(i)$ and $N_t(i)$, we obtain

by the tower property of expectation¹ that

$$\begin{aligned}
& \mathbb{E} \left[\exp \left(\iota \lambda_1 \sum_{\tau=1}^t \mathbf{1} \{A_\tau = i\} Y_\tau(i) + \iota \lambda_2 N_t(i) \right) \right] \\
&= \mathbb{E} \left[\prod_{\tau=1}^t \mathbb{E} [\exp (\iota \lambda_1 \mathbf{1} \{A_\tau = i\} Y_\tau(i) + \iota \mathbf{1} \{A_\tau = i\}) \mid \mathcal{H}_{\tau-1}] \right] \\
&= \mathbb{E} \left[\prod_{\tau=1}^t \left(\mathbf{1} \{A_\tau = i\} e^{\iota \lambda_2} \mathbb{E} [\exp(\iota \lambda_1 Y_\tau(i)) \mid \mathcal{H}_{\tau-1}] + \mathbf{1} \{A_\tau \neq i\} \right) \right]
\end{aligned}$$

because $N_t(i) = \sum_{\tau=1}^t \mathbf{1} \{A_\tau = i\}$ and $A_\tau \in \mathcal{H}_{\tau-1}$, that is, it is a function of the history. Then we observe that conditional on $\mathcal{H}_{\tau-1}$, $Y_\tau(i)$ and $Y'_\tau(i)$ have identical distribution, so that

$$\mathbb{E}[\exp(\iota \lambda_1 Y_\tau(i)) \mid \mathcal{H}_{\tau-1}] = \varphi_i(\lambda_1) = \mathbb{E}[\exp(\iota \lambda_1 Y'_\tau(i)) \mid \mathcal{H}_{\tau-1}].$$

Unrolling the above product, we therefore obtain

$$\begin{aligned}
& \mathbb{E} \left[\exp \left(\iota \lambda_1 \sum_{\tau=1}^t \mathbf{1} \{A_\tau = i\} Y_\tau(i) + \iota \lambda_2 N_t(i) \right) \right] \\
&= \mathbb{E} \left[\prod_{\tau=1}^t \left(\mathbf{1} \{A_\tau = i\} e^{\iota \lambda_2} \mathbb{E} [\exp(\iota \lambda_1 Y'_\tau(i)) \mid \mathcal{H}_{\tau-1}] + \mathbf{1} \{A_\tau \neq i\} \right) \right] \\
&= \mathbb{E} \left[\prod_{\tau=1}^t \mathbb{E} [\exp(\iota \lambda_1 Y'_\tau(i) \mathbf{1} \{A_\tau = i\} + \iota \lambda_2 \mathbf{1} \{A_\tau = i\}) \mid \mathcal{H}_{\tau-1}] \right] \\
&= \mathbb{E} \left[\exp \left(\iota \lambda_1 \sum_{\tau=1}^t \mathbf{1} \{A_\tau = i\} Y'_\tau(i) + \iota \mathbf{1} \{A_\tau = i\} \right) \right]
\end{aligned}$$

by again using the tower property of expectations. This gives the equality (18.2.1).

18.7 Further notes and references

Bandit problems, with their attendant need for analyses and procedures that consider both exploration and exploitation, have existed since at least 1933, beginning with Thompson's paper [179]. In the statistics literature, problems around the “design of experiments,” which resemble the linear bandit problems we outline in Section 18.3.4, consider choosing covariate vectors in linear regression to best estimate a regression vector θ [156], which Robbins considered in sequential settings in the 1950s [161]. The “modern” analysis of bandit problems took off with Lai and Robbins [130], who showed how confidence-based algorithms could achieve regret of order $K \log n / \Delta^2$, where $\Delta = \min_{i \neq i^*} \mu_{i^*} - \mu_i$, with attendant matching lower bounds.

More recent work in machine learning (of which there are far too many references to list) has been motivated by online auctions, medical scenarios, and, more recently, reinforcement learning problems. The books by Cesa-Bianchi and Lugosi [50], Bubeck and Cesa-Bianchi [42], and Lattimore and Szepesvári [132] provide excellent references, and several of our proofs follow those of Bubeck and Cesa-Bianchi [42]. Auer et al. [13] introduced the upper confidence bound (UCB)

¹That if X_t is a function of $\mathcal{H}_{t-1} \subset \mathcal{H}_t$, then $\mathbb{E}[X_1 \cdots X_N] = \mathbb{E}[\prod_{t=1}^N \mathbb{E}[X_t \mid \mathcal{H}_{t-1}]]$

algorithms and provided finite sample bounds on their regret, while Auer et al. [12] introduced EXP3 and its analysis. Interest in Thompson sampling reignited when Chapelle and Li [51] showed that, in spite of its (at the time) heuristic nature, it was competitive with established procedures, and Kaufmann et al. [124] provided the first analysis of the procedure. There are many connections between the particular cases we consider here and the broader field of reinforcement learning [175]; information-based perspectives remain an active area of research.

Our approach in Section 18.3 to Bayesian bandits follows the approach that Russo and Van Roy [164, 165, 166] pioneer. These analyses require a certain well-specification of the procedures, so that the prior and posterior distributions on the parameter are accurate, as otherwise the regret bounds fail to hold; extending Thompson sampling to apply even when the prior is unknown requires additional techniques (e.g. [3] or [132, Chapter 36]). That most analyses of bandit problems repose on some type of well-specification assumption—e.g., for linear bandits, that rewards $Y(a) = \langle \theta, a \rangle + \varepsilon$ —poses a problem for saying anything rigorous about their behavior in real-world scenarios. Researchers have begun to investigate this issue, showing that in the case of mis-specified models, finding the best approximation in the class can take time exponential in the dimension d of the approximator [69], while slightly weaker approximation guarantees admit efficient algorithms [133, 93]. Many questions in this direction remain open.

The approach we take in Section 18.5 derives from a variety of sources. At a high level, Section 18.5.1 is a distilled and slightly weaker version of Foster et al.’s *decision estimation coefficient* [91, 92]. We, however, take a perspective a little closer to the modulus of continuity of a statistic with respect to Hellinger distance, as in Chapter 13. The idea of using the empirical distribution of plays A_t to lower bound the regret as n -times the minimax risk (18.5.2) we learned from Bubeck and Cesa-Bianchi [42, Ch. 3.3], whose minimax approach also bears similarities to that we use in proving Theorem 18.5.1. The results using Assouad’s method to lower bound the expected regret of action A_{n+1} —the pure exploration case—appear to be new, though they take inspiration from Lattimore and Szepesvári [132, Chapter 24], who focus on lower bounds on the regret for both ℓ_∞ and ℓ_2 -style action sets; the technique with Assouad’s method also yields sharper constants.

JCD Comment: Also say something about causal inference!

JCD Comment: Can we do something like a treatment with covariates or something? Two actions, $a \in \{0, 1\}$. Get feedback $(X, Y(a))$, $X \in \mathbb{R}^d$. Define $\phi(a, x) = (a, x)$, and

$$(\hat{\tau}, \hat{\theta}) = \operatorname{argmin} P_n(\langle \phi(A, X), (\tau, \theta) \rangle - Y)^2$$

then for $H = P_n \phi \phi^\top$ and $\varepsilon = (Y - \langle \phi, (\tau, \theta) \rangle)$ we have

$$\sqrt{n}(\hat{\tau} - \tau) \overset{d}{\rightsquigarrow} \mathbf{N}\left(0, e_1^\top H^{-1} P \varepsilon^2 \phi \phi^\top H^{-1} e_1\right)$$

and naive estimator

$$\hat{\tau}_{\text{naive}} = \frac{2}{n} \sum_{i: A_i=1} Y_i - \frac{2}{n} \sum_{i: A_i=0} Y_i$$

has

$$\sqrt{n}(\hat{\tau}_{\text{naive}} - \tau) \overset{d}{\rightsquigarrow} \mathbf{N}(0, 4(\operatorname{Var}(Y(1)) + \operatorname{Var}(Y(0))))$$

18.8 Exercises

Exercise 18.1 (Convexity of quadratic-over-linear functions):

- (a) For $x, y \in \mathbb{R}$, define $h(x, y) = \frac{x^2}{y}$. Show that h is convex over the domain $y > 0$.
- (b) Let $x \in \mathbb{R}^n$ and $y \in \mathbb{R}$. Show that $h(x, y) = \|x\|^2 / y$ is convex over the domain $y > 0$ for any norm $\|\cdot\|$.
- (c) Let X and Y be vector spaces, and let $f : X \rightarrow \mathbb{R}^n$ and $g : Y \rightarrow \mathbb{R}$ be linear functions. Show that $h(x, y) = \|f(x)\|^2 / g(y)$ is convex over the domain $\{y \in Y \mid g(y) > 0\}$.
- (d) Assume f and g above are continuous linear functions, and that $f(x) \neq 0$ for all x . Show that h above is a closed convex function.

Exercise 18.2: Let $\sigma^2 > 0$ and $\tau^2 > 0$, $x_1, \dots, x_n \in \mathbb{R}^d$, and $y_1, \dots, y_n \in \mathbb{R}$. Let $\theta \in \mathbb{R}^d$ have density

$$\pi(\theta) \propto \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \langle x_i, \theta \rangle)^2 - \frac{1}{2\tau^2} \|\theta\|_2^2 \right).$$

Show that $\theta \sim \mathcal{N}(\hat{\theta}, \Sigma)$ for $\Sigma = \frac{1}{n}(\frac{1}{\tau^2}I_d + \frac{1}{n\sigma^2} \sum_{i=1}^n x_i x_i^\top)$ and

$$\hat{\theta} = \frac{1}{\sigma^2} \Sigma^{-1} \frac{1}{n} \sum_{i=1}^n x_i y_i = \operatorname{argmin}_{\theta} \left\{ \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \langle x_i, \theta \rangle)^2 + \frac{1}{2\tau^2} \|\theta\|_2^2 \right\}.$$

Exercise 18.3: We describe an approach using stochastic gradient methods to solve pure-exploration linear bandit problems. Assume that $\theta \in \mathbb{B}_2^d$ and let the action set $\mathcal{A} = \mathbb{B}_2^d$, and assume that $y(x) = \langle \theta, x \rangle + \varepsilon$ for ε an independent mean-zero variable with $\mathbb{E}[\varepsilon^2] \leq \sigma^2$. We wish to find a (potentially random) action $A \in \mathcal{A}$ so that $\mathbb{E}[\langle A, \theta \rangle] \leq -\|\theta\|_2 + o(1)$, that is, minimizing $\langle a, \theta \rangle$. Consider the following procedure: at iteration t , draw $w_t \sim \text{Uniform}(\mathbb{S}^{d-1})$, observe $y_t = \langle \theta, w_t \rangle + \varepsilon_t$, and then set the estimated subgradient $g_t = d \cdot w_t y_t$.

- (a) Show that $\mathbb{E}[g_t \mid \mathcal{H}_{t-1}] = \theta$ and $\mathbb{E}[\|g_t\|_2^2 \mid \mathcal{H}_{t-1}] \leq d\|\theta\|_2^2 + d^2\sigma^2$.
- (b) Let $\eta > 0$ be a stepsize (which you will determine). Define the iterates

$$\theta_t = \operatorname{argmin}_{\theta \in \mathbb{B}_2^d} \left\{ \langle g_t, \theta \rangle + \frac{1}{2\eta} \|\theta - \theta_{t-1}\|_2^2 \right\}$$

and $\bar{\theta}_n = \frac{1}{n} \sum_{t=1}^n \theta_t$. Give a setting of the stepsize η so that

$$\mathbb{E}[\langle \theta, \bar{\theta}_n \rangle] \leq -\|\theta\|_2 + O(1) \sqrt{\frac{d^2\sigma^2 + d}{n}}.$$

Exercise 18.4 (Exploration in linear bandits without noise): Consider a variant of the linear bandit problem with convex action set $\mathcal{A} \subset \mathbb{R}^d$ where at time t , given action x , we observe $y(x) = \langle \theta_t, x \rangle$ for θ_t chosen i.i.d. but satisfying $\mathbb{E}[\theta_t] = \theta$. Assume the action set $\mathcal{A} \supset \alpha\{-1, 1\}^d$ for some $\alpha > 0$ and that \mathcal{A} is convex.

(a) Show that the vector $g_t = \alpha^{-2} \langle \theta_t, A \rangle A$ for $A \sim \text{Uniform}(\alpha\{-1, 1\}^d)$ satisfies

$$\mathbb{E}[g_t] = \theta \quad \text{and} \quad \mathbb{E}[\|g_t\|_2^2] = d \cdot \mathbb{E}[\|\theta_t\|_2^2].$$

(b) Assume that the noise in θ_t is such that $\mathbb{E}[\|\theta_t\|_2^2] \leq \kappa^2$. Using the result of Exercise 18.3, show how to construct an action \hat{a}_n based on n observations such that

$$\mathbb{E}[\langle \theta, \hat{a}_n \rangle] \leq \inf_{a \in \mathcal{A}} \langle \theta, a \rangle + O(1)\kappa \text{diam}(\mathcal{A}) \sqrt{\frac{d}{n}}.$$

(c) Compare the result of part (b) to achievable convergence guarantees in the full information case, where we may observe θ_t , and the upper and lower bounds for cases in which $y(x) = \langle \theta, x \rangle + \varepsilon$.

Exercise 18.5: Show that information directed sampling without noise and sparse θ is really easy (better than UCB, say).

JCD Comment: Put in a few further examples, or leave as exercises?

- i. Linear bandits
- ii. Full information (or something similar)
- iii. We can pretty easily show that in the two-armed bandit case, a bound of $1/\Delta^2$ (where Δ is the gap in arms) is how many pulls are needed for each arm. Can also show that probability of error at time $n + 1$ is at least $\exp(-n\Delta^2)$ or so (Bretagnolle-Huber).

JCD Comment: We can do one on heavy-tailed (or at least just lighter MGFs)-based UCB algorithms.

JCD Comment: One on causality?

JCD Comment: Batched bandits?

Exercise 18.6: Prove a high-probability regret bound for UCB.

JCD Comment: Some experiments with Gaussians perhaps, and a mis-specified prior I think. Like what happens if the prior is too diffuse or not diffuse enough?

Chapter 19

Minimax games and Bayesian estimation

This final chapter explores a few of the many connections between online learning, probabilistic prediction, and Bayesian statistics, where we link the areas via information theoretic analyses. Within information theory, these problems classically arise from universal prediction, where one wishes to encode a sequence of random variables arriving sequentially from an unknown distribution nearly as well as if one knew the distribution. We use some of these ideas as motivation, but we will focus more on the statistical sides of the problem and connections with the theory of proper losses for prediction we develop in Chapter 14, referring to the bibliographic section for further exploration.

Consider the following minimax probabilistic game between nature, who chooses some distribution to sample from a set \mathcal{X} (for now finite), and a decision maker or player, who wishes to model the sampling distribution as well as possible. We, the player, first choose a distribution Q on \mathcal{X} with probability mass function q , and then nature chooses a $P \in \mathcal{P}$, where \mathcal{P} is a collection of distributions on \mathcal{X} . Nature draws $X \sim P$, and upon revealing a realization x , we suffer log loss $-\log q(x)$ and expected loss $\mathbb{E}_P[-\log q(X)]$. In our game we thus suffer (worst-case) expected loss

$$\sup_{P \in \mathcal{P}} \mathbb{E}_P[-\log q(X)] = \sup_{P \in \mathcal{P}} \sum_{x \in \mathcal{X}} p(x) \log \frac{1}{q(x)}, \quad (19.0.1)$$

and to play optimally, we wish to solve the minimax problem

$$\underset{Q}{\text{minimize}} \sup_{P \in \mathcal{P}} \mathbb{E}_P[-\log q(X)]. \quad (19.0.2)$$

Instead of using the objective (19.0.2), it is frequently sensible to relativize losses to \mathcal{P} : we wish to compete against the best $P \in \mathcal{P}$, and (as in Chapter 17) seek to play distributions Q where we have little *regret* relative to the losses we would suffer if we had played the optimal P . Then instead of the problem (19.0.2), we consider the objective

$$\sup_{P \in \mathcal{P}} \mathbb{E}_P[-\log q(X)] - \inf_{Q \in \mathcal{P}} \mathbb{E}_P[-\log q(X)] = \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\log \frac{p(X)}{q(X)} \right],$$

the KL-divergence between P and Q . More abstractly, we can revisit the prediction setting of Chapter 14, where we for a proper loss ℓ , we suffer $\ell(Q, x)$ for predicting distribution Q on realization

x ; the problem (19.0.2) corresponds to the logarithmic loss. Then we wish to minimize the worst-case expected regret of playing Q instead of the optimal P , that is, solve the game

$$\text{minimize}_Q \sup_{P \in \mathcal{P}} \{\mathbb{E}_P[\ell(Q, X)] - \mathbb{E}_P[\ell(P, X)]\}. \quad (19.0.3)$$

Such minimax games are a focus in the game theory, economics, and optimization literatures, and as we shall see, they bring insights to information theory and statistics.

We can also imagine a variant of the game (19.0.3). To make the connections with parameter estimation and Bayesian statistics notationally clearer, let the collection $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$ be parameterized by θ (this is no loss of generality). Then instead of choosing a distribution Q and allowing nature to choose a distribution P_θ , we could switch the order of the game: nature first chooses prior distribution π on θ , and without seeing θ (but with knowledge of the distribution π) we choose the predictive distribution Q . This leads to the *Bayesian regret*, which is simply the expectation

$$\int_{\Theta} (\mathbb{E}_{P_\theta}[\ell(Q, X)] - \mathbb{E}_{P_\theta}[\ell(P_\theta, X)]) \pi(\theta) d\theta.$$

Let P_π denote the marginal distribution over X obtained by drawing $\theta \sim \pi$ and then drawing $X \sim P_\theta$. When the loss ℓ is proper (recall Chapter 14.1), P_π minimizes the Bayesian regret. Nature, acting adversarially, may choose a prior to make this regret as large as possible, but the Q player always plays with knowledge of nature's choice. This leads to the main question of this chapter: when are the (worst case) Bayesian regret and the game (19.0.3) where Q plays first equivalent?

This chapter investigates the problems (19.0.2) and (19.0.3) and their variants. We begin by motivating the problems in two ways: first, in Section 19.1, as a form of robust Bayesian procedure, which leads to classical maximum entropy estimators. In Section 19.2, we overview some results on universal and sequential prediction, connecting coding problems to the games (19.0.2) and (19.0.3). We then present a fundamental duality result: that (for many classes \mathcal{P} of distributions), it does not matter which player goes first in problem (19.0.3), and similarly in problem (19.0.2): there are robust choices of Q that nature—the P player—can, essentially, gain no advantage against. For the remainder of the chapter, we then provide analyses of Bayesian statistical procedures, showing how they provide asymptotic guarantees on the log loss $-\mathbb{E}_P[\log q(X_1^n)]$ when $X_i \stackrel{\text{iid}}{\sim} P$, giving explicit ways to (asymptotically) find such optimal Q distributions.

19.1 Robust Bayesian procedures and maximum entropy

Consider the setting of Bayesian statistics, where a data analyst has a prior π on the states of nature, and wishes to perform some type of inference to estimate the “true” state of nature using this prior. For simplicity, let us assume $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$, where π is some density on Θ . Then the Bayesian chooses Q to maximize the log-likelihood of observed data—to best approximate P_θ —on average over the possible states θ of nature, that is, to minimize

$$\int \mathbb{E}_{P_\theta}[-\log q(X)] \pi(\theta) d\theta.$$

A “robust” Bayesian chooses the worst possible prior π on Θ from a collection of potential priors, yielding the robust loss

$$\sup_{\pi \in \Pi} \int \mathbb{E}_{P_\theta}[-\log q(X)] \pi(\theta) d\theta.$$

Abstracting away the particular representation to simply consider the family $\mathcal{P} = \{P_\pi\}_{\pi \in \Pi}$ of distributions, where we recall the mixture $P_\pi(A) = \int P_\theta(A)\pi(\theta)d\theta$, thus yields the worst case log loss (19.0.1). So we think of problem (19.0.2) as seeking *robust Bayes* acts against \mathcal{P} : we wish to solve

$$\text{minimize}_Q \sup_{\pi \in \Pi} \int \mathbb{E}_{P_\theta}[-\log q(X)]\pi(\theta)d\theta. \quad (19.1.1)$$

We say that Q is a robust Bayes procedure for the prior distributions Π if it minimizes the supremum risk (19.1.1); that is, it achieves the best possible logarithmic loss measured in a uniform sense against all distributions $P \in \mathcal{P} = \{P_\pi\}_{\pi \in \Pi}$.

Consider temporarily a discrete set \mathcal{X} , and let $H(P) = -\sum_x p(x) \log p(x)$ be the usual (Shannon) entropy. Because the KL-divergence is always nonnegative, it is immediately apparent that for any distribution P ,

$$\mathbb{E}_P[-\log q(X)] = \mathbb{E}_P[-\log p(X)] + D_{\text{kl}}(P\|Q) \geq H(P),$$

and equality holds if and only if $Q = P$. Thus, abstracting away the particular prior family to consider collections \mathcal{P} of distributions, if we switch the order of play in the game (19.1.1), forcing nature to play first, then

$$\sup_{P \in \mathcal{P}} \inf_Q \mathbb{E}_P[-\log q(X)] = \sup_{P \in \mathcal{P}} H(P),$$

the maximum entropy in \mathcal{P} . Then our leading question in the introductory remarks of the chapter—when does it not matter whether the Q player (statistician) or P player (nature) plays first in the game (19.1.1)?—becomes one of duality: when do we have

$$\sup_{P \in \mathcal{P}} \inf_Q \mathbb{E}_P[-\log q(X)] = \inf_Q \sup_{P \in \mathcal{P}} \mathbb{E}_P[-\log q(X)]? \quad (19.1.2)$$

19.1.1 A digression on min-max games

We provide a brief digression on min-max dualities before continuing. For any function $L : U \times V \rightarrow \mathbb{R}$, where U and V are arbitrary spaces, the min-max game considers solving

$$\text{minimize}_{u \in U} \sup_{v \in V} L(u, v) \quad \text{and} \quad \text{maximize}_{v \in V} \inf_{u \in U} L(u, v), \quad (19.1.3)$$

where we think of a u -player and a v -player choosing u (respectively v), after which the other player chooses a best-response. The game has a value if the order of play does not matter, that is, when $\inf_u \sup_v L(u, v) = \sup_v \inf_u L(u, v)$. In general, no matter U , V , and L , we always have the *weak min-max inequality*

$$\sup_{v \in V} \inf_{u \in U} L(u, v) \leq \inf_{u \in U} \sup_{v \in V} L(u, v). \quad (19.1.4)$$

Indeed, for any $u_0 \in U$ and $v_0 \in V$, we certainly have $\inf_{u \in U} L(u, v_0) \leq \sup_{v \in V} L(u_0, v)$, and taking the supremum over v_0 on the left and infimum over u_0 on the right implies inequality (19.1.4). A pair $(u^*, v^*) \in U \times V$ is a *saddle point* for the game if

$$\sup_{v \in V} L(u^*, v) \leq L(u^*, v^*) \leq \inf_{u \in U} L(u, v^*), \quad (19.1.5)$$

the saddle nomenclature following as (u^*, v^*) simultaneously minimizes and maximizes L . For such a point, we necessarily have

$$\inf_{u \in U} \sup_{v \in V} L(u, v) \leq \sup_{v \in V} L(u^*, v) \leq L(u^*, v^*) \leq \inf_{u \in U} L(u, v^*) \leq \sup_{v \in V} \inf_{u \in U} L(u, v),$$

so in view of the weak min-max inequality (19.1.4), each inequality is necessarily an equality above. We capture this as a proposition for later reference:

Proposition 19.1.1. *If a point (u^*, v^*) is a saddle point for the game (19.1.3), then the game has value $L(u^*, v^*)$, independent of the order of play, and*

$$L(u^*, v^*) = \inf_{u \in U} \sup_{v \in V} L(u, v) = \sup_{v \in V} \inf_{u \in U} L(u, v).$$

Moreover, u^* minimizes $\sup_{v \in V} L(u, v)$ over $u \in U$, and v^* maximizes $\inf_{u \in U} L(u, v)$ over $v \in V$.

Classical results in convex optimization and game theory provide sufficient conditions for such saddle points to exist and, more generally, for equality to hold in inequality (19.1.4). Exercise 17.3 shows how to use online convex optimization techniques to prove the classical von Neumann minimax theorem, that when $U \subset \mathbb{R}^m$ and $V \subset \mathbb{R}^n$ are compact convex sets and $L(u, v) = \langle u, Av \rangle$ for some matrix A , then there indeed exists a saddle point. Appendix C.4 reviews the main results in (finite-dimensional) convex analysis on existence of saddle points, though because of the particular structure of the regret minimization game (19.0.3), we shall be able to give a more direct “information-theoretic” proof of duality, which will have the advantage that it applies equally to finite and infinite-dimensional problems.

19.1.2 Saddle points for maximum entropy

Returning to the log-loss game (19.0.2) and robust Bayesian estimation problem (19.1.1), under appropriate conditions, we expect to have a saddle point, so that there are Q^*, P^* such that

$$\sup_{P \in \mathcal{P}} \mathbb{E}_P[-\log q^*(X)] \leq \mathbb{E}_{P^*}[-\log q^*(X)] \leq \inf_Q \mathbb{E}_{P^*}[-\log q(X)],$$

which immediately implies the equality (19.1.2) by Proposition 19.1.1. In this case, when we do indeed have a saddle point, the saddle point thus necessarily achieves the maximum entropy $H(P^*) = \sup_{P \in \mathcal{P}} H(P)$, which is thus the value for the game (19.0.2). By the defining criterion for a saddle point, q^* evidently uniquely minimizes $\mathbb{E}_{P^*}[-\log q(X)] = H(P^*) + D_{\text{kl}}(P^* \| Q)$, so we see that Q^* maximizes the entropy over \mathcal{P} as well, so that maximum entropy becomes the robust Bayes act against \mathcal{P} : it solves the worst-case problem (19.0.2) and maximizes entropy. In summary, for the log-loss game, any saddle point takes the form

$$(P^*, P^*),$$

and the distribution P^* maximizes the entropy.

In Section 19.6, we shall show that saddle points typically exist for the expected regret game (19.0.3) when the losses ℓ are strictly proper, connecting the development of our loss functions in Chapter 14 to the game formulations (19.0.2) and (19.0.3). As corollaries, we shall develop the existence of saddle points in the log-loss game and more general scenarios. First, however, we show that exponential family models provide a clean treatment of the log-loss game without any sweat.

19.1.3 Exponential family models as robust Bayesian procedures

We now show how the minimax game (19.0.1) naturally gives rise to exponential family models, so that our now familiar exponential family distributions from Chapter 3 are in fact robust Bayes

procedures against certain families \mathcal{P} of distributions. To that end, let $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$ be a statistic, and for $\alpha \in \mathbb{R}^d$ consider the mean-value constrained set

$$\mathcal{P} = \mathcal{P}_\alpha^{\text{lin}} := \{\text{distributions } P \text{ on } \mathcal{X} \text{ s.t. } \mathbb{E}_P[\phi(X)] = \alpha\}.$$

Continuing to consider the case that \mathcal{X} is discrete, let P_θ denotes the exponential family distribution with p.m.f. $p_\theta(x) = h(x) \exp(\langle \theta, \phi(x) \rangle - A(\theta))$, where h denotes a carrier. We have the following.

Proposition 19.1.2. *Let the conditions above hold. If $\mathbb{E}_{P_\theta}[\phi(X)] = \alpha$, then*

$$\inf_Q \sup_{P \in \mathcal{P}_\alpha^{\text{lin}}} \mathbb{E}_P[-\log q(X)] = \sup_{P \in \mathcal{P}_\alpha^{\text{lin}}} \mathbb{E}_P[-\log p_\theta(X)] = \sup_{P \in \mathcal{P}_\alpha^{\text{lin}}} \inf_Q \mathbb{E}_P[-\log q(X)].$$

Proof First, note that

$$\begin{aligned} \sup_{P \in \mathcal{P}_\alpha^{\text{lin}}} \mathbb{E}_P[-\log p_\theta(X)] &= \sup_{P \in \mathcal{P}_\alpha^{\text{lin}}} \mathbb{E}_P[-\langle \phi(X), \theta \rangle + A(\theta)] \\ &= -\langle \alpha, \theta \rangle + A(\theta) = \mathbb{E}_{P_\theta}[-\langle \theta, \phi(X) \rangle + A(\theta)] = H(P_\theta), \end{aligned}$$

where H denotes the Shannon entropy, for any distribution $P \in \mathcal{P}_\alpha^{\text{lin}}$. Moreover, for any $Q \neq P_\theta$, we have

$$\sup_P \mathbb{E}_P[-\log q(X)] \geq \mathbb{E}_{P_\theta}[-\log q(X)] > \mathbb{E}_{P_\theta}[-\log p_\theta(X)] = H(P_\theta),$$

where the inequality follows because $D_{\text{kl}}(P_\theta \| Q) > 0$. This shows the first equality in the proposition.

For the second equality, note that

$$\inf_Q \mathbb{E}_P[-\log q(X)] = \inf_Q \underbrace{\mathbb{E}_P \left[\log \frac{p(X)}{q(X)} \right]}_{=0} - \mathbb{E}_P[\log p(X)] = H(P).$$

But we know from our standard maximum entropy results (Theorem 14.4.7) that P_θ maximizes the entropy over $\mathcal{P}_\alpha^{\text{lin}}$, that is, $\sup_{P \in \mathcal{P}_\alpha^{\text{lin}}} H(P) = H(P_\theta)$. \square

The same result immediately holds beyond discrete sets \mathcal{X} , which we state here without proof, as it is more or less completely identical to the proof of Proposition 19.1.2. Let ν be some measure on \mathcal{X} (e.g., the Lebesgue measure $d\nu(x) = dx$), and let

$$\mathcal{P}_\alpha^{\text{lin}} := \{\text{distributions } P \ll \nu \text{ on } \mathcal{X} \text{ s.t. } \mathbb{E}_P[\phi(X)] = \alpha\}.$$

Then the following generalization of Proposition 19.1.2 holds.

Proposition 19.1.3. *Let $\mathcal{E} = \{P_\theta\}_{\theta \in \mathbb{R}^d}$ denote the exponential family with densities $p_\theta(x) = \exp(\langle \theta, \phi(x) \rangle - A(\theta))$ with respect to ν . Then if $\mathbb{E}_{P_\theta}[\phi(X)] = \alpha$, P_θ is the uniquely robust Bayes distribution against $\mathcal{P}_\alpha^{\text{lin}}$.*

Exercise 19.5 asks you to show that the proof of Proposition 19.1.2 extends straightforwardly.

In short: maximum entropy is equivalent to robust prediction procedures for linear families of distributions $\mathcal{P}_\alpha^{\text{lin}}$. In turn, this is equivalent to maximum likelihood estimation in exponential families, which corresponds to moment-matching (3.2.3).

19.2 The coding game and sequential prediction

An additional motivation for the log-loss game (19.0.2) comes from coding problems. In this case, we assume we receive n random variables $X_i \stackrel{\text{iid}}{\sim} P$, where P is unknown, and suffer the sequential prediction loss

$$\mathbb{E}_P[-\log q(X_1^n)] = \sum_{i=1}^n \mathbb{E}_P \left[\log \frac{1}{q(X_i | X_1^{i-1})} \right],$$

which corresponds to predicting X_i given X_1^{i-1} as well as possible, even when the X_i follow an (unknown or adversarially chosen) distribution P . The connection between instantaneously decodable codes and probability distributions underlying our original interpretations of the entropy, as in Chapter 2.4, provides a concrete grounding for the minimax game (19.0.2).

Example 19.2.1 (The coding game): Let the set \mathcal{X} be finite or countable, and consider the problem of encoding \mathcal{X} into $\{0, 1\}$ -valued sequences using as few bits as possible. In this case, the Kraft inequality (recall Theorem 2.4.2) tells us that if $C : \mathcal{X} \rightarrow \{0, 1\}^*$ is a uniquely decodable code, and $\ell_C(x)$ denotes the length of the encoding for the symbol $x \in \mathcal{X}$, then

$$\sum_x 2^{-\ell_C(x)} \leq 1.$$

Thus, we may define the distribution $q_C(x) = 2^{-\ell_C(x)} / \sum_x 2^{-\ell_C(x)}$, and for any sequence x_1^n , we have

$$-\log_2 q_C(x_1^n) = \sum_{i=1}^n \left[\ell_C(x_i) + \log \sum_x 2^{-\ell_C(x)} \right] \leq \sum_{i=1}^n \ell_C(x_i).$$

Conversely, given any distribution q on $x_1^n \in \mathcal{X}^n$, the function $\ell : \mathcal{X}^n \rightarrow \mathbb{N}$ with $\ell(x_1^n) := \lceil -\log_2 q(x_1^n) \rceil$ evidently satisfies $\sum_x 2^{-\ell(x_1^n)} \leq \sum_x 2^{\log_2 q(x_1^n)} = 1$, and so there exists a binary prefix code $C : \mathcal{X}^n \rightarrow \{0, 1\}^*$ with the length function $\ell_C(x_1^n) = \lceil -\log_2 q(x_1^n) \rceil$ by the converse part of the Kraft inequality in Theorem 2.4.2. Thus for $X_1^n \sim P$,

$$\mathbb{E}_P[\ell_C(X_1^n)] = \mathbb{E}_P[\lceil -\log_2 q(X_1^n) \rceil] \leq 1 + \mathbb{E}_P[-\log_2 q(X_1^n)].$$

The minimax game (19.0.2) thus corresponds to a coding game where we attempt to choose a distribution Q (or sequential coding scheme C) that has as small an expected length as possible, uniformly over distributions P . \diamond

When we have sequences x_1, \dots, x_n , the minimax log-loss game (19.0.2) thus is equivalent to a sequential setting in which we attempt to predict symbols, or equivalently, encode the sequence x_1^n online. Rather than measuring performance simply by the losses $-\log q(x)$, however, we measure against reference distributions P , and ask whether, we can predict the sequence (nearly) as well as if we knew the true distribution of the data. Or, in more general settings, we would like to predict the data as well as all predictive distributions P from some family of distributions \mathcal{P} , even if *a priori* we know little about the coming sequence of data.

This leads us to an online setting that parallels the online convex optimization scenarios we investigated in Chapter 17. With the log loss $\ell(q, x) = -\log q(x)$, we even are in the setting of that chapter, but the explosion of $-\log q(x)$ near the boundaries of the probability simplex necessitates a different set of tools. We consider two versions of the sequential prediction game: adversarial and probabilistic. For both of the following definitions of sequential prediction games, we assume that p and q are densities or probability mass functions in the case that \mathcal{X} is continuous or discrete (this is no real loss of generality) for distributions P and Q .

Adversarial regret We begin with the adversarial case. Given a sequence $x_1^n \in \mathcal{X}^n$, the *regret* of the distribution Q for the sequence x_1^n with respect to the distribution P is

$$\text{Reg}(Q, P, x_1^n) := \log \frac{1}{q(x_1^n)} - \log \frac{1}{p(x_1^n)} = \sum_{i=1}^n \log \frac{1}{q(x_i | x_1^{i-1})} - \log \frac{1}{p(x_i | x_1^{i-1})}, \quad (19.2.1)$$

where we have written it as the sum over $q(x_i | x_1^{i-1})$ to emphasize the sequential nature of the game. The quantity (19.2.1) measures how much we “regret” playing a distribution Q over the alternate distribution P . In the context of the coding game in Example 19.2.1, this corresponds to the realized number of bits necessary to encode x_1^n over a putative encoding through p . Associated with the regret of the sequence x_1^n is the *adversarial regret* of Q with respect to the family \mathcal{P} of distributions, which we typically call the regret as in Chapter 17, which is

$$\mathfrak{R}_n^{\mathcal{X}}(Q, \mathcal{P}) := \sup_{P \in \mathcal{P}, x_1^n \in \mathcal{X}^n} \text{Reg}(Q, P, x_1^n). \quad (19.2.2)$$

Redundancy A less adversarial problem is to minimize the *expected regret* with respect to a distribution P , which for the log loss gets the special name *redundancy*. We thus define

$$\text{Red}_n(Q, P) := \mathbb{E}_P \left[\log \frac{1}{q(X_1^n)} - \log \frac{1}{p(X_1^n)} \right] = D_{\text{kl}}(P \| Q), \quad (19.2.3)$$

where the dependence on n is implicit in the KL-divergence. The name choice follows from the connections to coding games in Example 19.2.1:

Example 19.2.2 (Example 19.2.1 on coding, continued): For any p.m.f.s p and q on the set \mathcal{X} , we can define coding schemes C_p and C_q with code lengths

$$\ell_{C_p}(x) = \left\lceil \log \frac{1}{p(x)} \right\rceil \quad \text{and} \quad \ell_{C_q}(x) = \left\lceil \log \frac{1}{q(x)} \right\rceil.$$

Conversely, given (uniquely decodable) encoding schemes C_p and $C_q : \mathcal{X} \rightarrow \{0, 1\}^*$, the functions $p_{C_p}(x) = 2^{-\ell_{C_p}(x)}$ and $q_{C_q}(x) = 2^{-\ell_{C_q}(x)}$ satisfy $\sum_x p_{C_p}(x) \leq 1$ and $\sum_x q_{C_q}(x) \leq 1$ by the Kraft-McMillan inequality (Theorem 2.4.2). Thus, the redundancy of Q with respect to P is the additional number of bits required to encode variables distributed according to P when we assume they have distribution Q , that is, how redundant our implicit encoding by Q is:

$$\begin{aligned} \text{Red}_n(Q, P) &= \sum_{i=1}^n \mathbb{E}_P \left[\log \frac{1}{q(X_i | X_1^{i-1})} - \log \frac{1}{p(X_i | X_1^{i-1})} \right] \\ &= \sum_{i=1}^n \mathbb{E}_P [\ell_{C_q}(X_i)] - \mathbb{E}_P [\ell_{C_p}(X_i)], \end{aligned}$$

where $\ell_C(x)$ denotes the number of bits C uses to encode x . As in Section 2.4.1, the code $\lceil -\log p(x) \rceil$ is (essentially) optimal. \diamond

As another example, we may consider a filtering or prediction problem for a linear system.

Example 19.2.3 (Prediction in a linear system): Suppose we believe that a sequence of random variables $X_i \in \mathbb{R}^d$ are Markovian, where X_i given X_{i-1} is normally distributed with mean $AX_{i-1} + g$, where A is an unknown matrix and $g \in \mathbb{R}^d$ is a constant drift term. Concretely, we assume $X_i \sim \mathcal{N}(AX_{i-1} + g, \sigma^2 I_{d \times d})$, where we assume σ^2 is fixed and known. For our class of predicting distributions Q , we may look at those that at iteration i predict $X_i \sim \mathcal{N}(\mu_i, \sigma^2 I)$. This yields regret

$$\text{Reg}(Q, P, x_1^n) = \sum_{i=1}^n \frac{1}{2\sigma^2} \|\mu_i - x_i\|_2^2 - \frac{1}{2\sigma^2} \|Ax_{i-1} + g - x_i\|_2^2,$$

while the redundancy is

$$\text{Red}_n(Q, P) = \frac{1}{2\sigma^2} \sum_{i=1}^n \mathbb{E}[\|AX_{i-1} + g - \mu_i(X_1^{i-1})\|_2^2],$$

assuming that P is the linear Gaussian Markov chain. \diamond

19.3 Expected regret, information capacity, and redundancy

That the log-loss admits a particularly elegant reformulation in terms of the KL-divergence is but a special case of more general results connecting divergences and the proper losses of Chapter 14, of which the log-loss is one. Thus, we consider minimizing worst-case expected regret as in problem (19.0.3). Now, for a general (proper) loss function ℓ , we define the regret

$$\text{Reg}(Q, P, x_1^n) := \sum_{i=1}^n \ell(Q(\cdot | x_1^{i-1}), x_i) - \ell(P(\cdot | x_1^{i-1}), x_i),$$

where $\ell(P, x_i)$ indicates loss suffered on the point x_i when the distribution P over X_i is played, and $P(\cdot | x_1^{i-1})$ denotes the conditional distribution of X_i given x_1^{i-1} according to P . The particular loss should clear from context in most of our discussion.

Taking expectations, we obtain the *expected regret*

$$\begin{aligned} \text{EReg}_n(Q, P) &:= \mathbb{E}_P \left[\sum_{i=1}^n \ell(Q(\cdot | X_1^{i-1}), X_i) - \ell(P(\cdot | X_1^{i-1}), X_i) \right] \\ &= \sum_{i=1}^n \mathbb{E}_P [\ell(Q(\cdot | X_1^{i-1}), X_i) - \ell(P(\cdot | X_1^{i-1}), X_i)], \end{aligned}$$

where the expectation is taken over $X_1^n \sim P$. Because we focus on proper losses, the expected regret is always nonnegative, and when ℓ is strictly proper, it is positive unless $Q = P$. The worst-case expected regret with respect to a class \mathcal{P} is then $\sup_{P \in \mathcal{P}} \text{EReg}_n(Q, P)$.

To simplify notation, let us consider the case that $n = 1$ (or, equivalently, that we simply measure losses $\ell(Q, X_1^n)$ against an entire vector X_1^n). Then minimizing the worst-case expected regret is identical to the problem (19.0.3), and

$$\text{EReg}(Q, P) = \mathbb{E}_P[\ell(Q, X)] - \mathbb{E}_P[\ell(P, X)],$$

because ℓ is proper. We can make an essentially complete analogy with the redundancy (19.2.3), that is, that the expected regret is a divergence between distributions. To do this, recall the Savage

representation of proper losses from Chapter 14.2. Corollaries 14.2.5 and 14.2.9 show the loss ℓ is (strictly) proper if and only if there exists a (strictly) convex Ω such that

$$\ell(Q, x) = -\Omega(Q) - \langle \nabla \Omega(Q), e_x - Q \rangle$$

(where we abuse notation slightly to write inner products), and for this negative entropy Ω ,

$$D_\Omega(P, Q) = \mathbb{E}_P[\ell(Q, X)] - \mathbb{E}_P[\ell(P, X)],$$

where we recall the Bregman divergence $D_\Omega(P, Q) = \Omega(P) - \Omega(Q) - \langle \nabla \Omega(Q), P - Q \rangle$. That is, the expected regret

$$\text{EReg}(Q, P) = D_\Omega(P, Q) \tag{19.3.1}$$

is the divergence between P and Q , as in the case (19.2.3) of the log loss with the KL-divergence, and the representation (19.3.1) via divergence will be the key to the dualities we show.

We consider the min-max game formulation, but set notation to evoke estimating parameters θ of a model of interest. Indeed, without loss of generality, we can identify $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$ for some space Θ (for example, taking Θ to be in one-to-one mapping with probability distributions). Then letting $\Pi(\Theta)$ be the collection of probability distributions on Θ , we evidently have

$$\sup_{P \in \mathcal{P}} \text{EReg}(Q, P) = \sup_{\pi \in \Pi(\Theta)} \int \text{EReg}(Q, P_\theta) d\pi(\theta),$$

and the rightmost quantity is linear in π . Then we can ask for a distribution Q minimizing the worst-case expected regret, and similarly ask when it does not matter whether we choose Q first and nature chooses π or vice versa: do we have the duality

$$\inf_Q \sup_\pi \int \text{EReg}(Q, P_\theta) d\pi(\theta) = \sup_\pi \inf_Q \int \text{EReg}(Q, P_\theta) d\pi(\theta)?$$

19.3.1 Information capacity and regret duality

Recall now Chapter 14.1.3, where we defined the *information in an experiment* as the generalized entropy reduction that observing a variable X provides for Y . Adapting our notation here, assume the Markov chain $T \rightarrow X$, where we draw $T \sim \pi$ for a (prior) probability distribution Θ , and conditional on $T = \theta$, we draw $X \sim P_\theta$. Then we defined the information in the experiment (14.1.7) by $I_\ell(T; X) = H_\ell(X) - H_\ell(X | T)$. To make apparent the importance of the prior distribution π on T , we incorporate it into our notation and define the information X carries about T according to the loss ℓ via

$$\begin{aligned} I_\ell(\pi; X) &:= \inf_Q \int \text{EReg}(Q, P_\theta) d\pi(\theta) \\ &= \inf_Q \int \mathbb{E}_{P_\theta}[\ell(Q, X)] d\pi(\theta) - \int \mathbb{E}_{P_\theta}[\ell(P_\theta, X)] d\pi(\theta) \\ &= H_\ell(X) - H_\ell(X | T), \end{aligned} \tag{19.3.2}$$

where the final equality uses the generalized entropy (14.1.5) and tacitly draws $T \sim \pi$.

Because the loss ℓ is proper, the infimum over Q in the definition (19.3.2) is always attained by the marginal distribution

$$P_\pi(A) := \int P_\theta(A) d\pi(\theta), \tag{19.3.3}$$

that is, P_π is the marginal distribution on X after drawing $T \sim \pi$ and then $X \sim P_T$. To see this, simply note that $\int \mathbb{E}_{P_\theta}[\ell(Q, X)]d\pi(\theta) = \mathbb{E}_{P_\pi}[\ell(Q, X)]$, which P_π minimizes by propriety. Recalling the representation of expected regret as the divergence (19.3.1) for the (generalized) negative entropy $\Omega(P) = -\inf_P \mathbb{E}_P[\ell(P, X)]$, we thus may represent the information as the average Bregman divergence

$$I_\ell(\pi; X) = \inf_Q \int D_\Omega(P_\theta, Q)d\pi(\theta) = \int D_\Omega(P_\theta, P_\pi)d\pi(\theta).$$

This should be familiar, as it generalizes the result (9.4.4) for the Shannon mutual information (corresponding to the log loss) that $I(T; X) = \int D_{\text{kl}}(P_\theta \| P_\pi) d\pi(\theta)$.

The notation (19.3.2) makes clear the dependence on the prior π , and so for a collection $\Pi \subset \Pi(\Theta)$ of probability distributions on Θ , we define the *capacity* of the channel $T \rightarrow X$ by

$$C_\ell(\Pi) := \sup_{\pi \in \Pi} I_\ell(\pi; X). \quad (19.3.4)$$

This measures the most information that the channel $T \rightarrow X$ can possibly contain when we know that T is drawn from the prior $\pi \in \Pi$. To rephrase our questions of duality as a saddle point problem as in Section 19.1.1, we define the expected loss gap (regret) with respect to the prior π ,

$$L(Q, \pi) := \int \mathbb{E}_{P_\theta}[\ell(Q, X)]d\pi(\theta) - \int \mathbb{E}_{P_\theta}[\ell(P_\theta, X)]d\pi(\theta) = \int D_\Omega(P_\theta, Q)d\pi(\theta). \quad (19.3.5)$$

Then the worst-case expected regret (redundancy) of a distribution Q on \mathcal{X} is

$$\sup_{\pi \in \Pi} L(Q, \pi) = \sup_{\pi \in \Pi} \int D_\Omega(P_\theta, Q)d\pi(\theta),$$

while the capacity (19.3.4) of the family $\{P_\theta\}_{\theta \in \Theta}$ is

$$C_\ell(\Pi) := \sup_{\pi \in \Pi} I_\ell(\pi; X) = \sup_{\pi \in \Pi} \inf_Q L(Q, \pi) = \sup_{\pi \in \Pi} L(P_\pi, \pi),$$

where we recall that P_π denotes the marginal distribution (19.3.3) and thus satisfies

$$\inf_Q L(Q, \pi) = L(P_\pi, \pi). \quad (19.3.6)$$

We can therefore rephrase our questions of duality as the following *regret/capacity duality* question: under what circumstances do we have min-max equality for the saddle objective (19.3.5),

$$\inf_Q \sup_{\pi \in \Pi} L(Q, \pi) = \sup_{\pi \in \Pi} \inf_Q L(Q, \pi) = \sup_{\pi \in \Pi} I_\ell(\pi; X)? \quad (19.3.7)$$

Because $\Pi \subset \Pi(\Theta)$ may be infinite dimensional, this equality does not follow so immediately from classical minimax theorems. Nonetheless, the duality (19.3.7) does indeed typically hold.

A key intermediate result is that when there is a prior π^* achieving the capacity, then we are guaranteed a saddle point, and so long as ℓ is strictly proper, we have uniqueness. Because we wish to actually achieve minimizers, we make a minor restriction (as we will see, this is in practice not important) to assume the losses admit a one-dimensional lower semicontinuity. So we say that $D_\Omega(P, Q)$ is *directionally lower semicontinuous in Q* if for any two distributions Q_0, Q_1 ,

$$\liminf_{\lambda \downarrow 0} D_\Omega(P, (1 - \lambda)Q_0 + \lambda Q_1) \geq D_\Omega(P, Q_0). \quad (19.3.8)$$

Under this condition, achieving the maximum capacity is sufficient to guarantee the duality (19.3.7), and even more, the existence of saddle points:

Lemma 19.3.1. *Let Π be a convex collection of distributions on Θ . Assume that the capacity $C_\ell(\Pi) := \sup_{\pi \in \Pi} I_\ell(\pi; X)$ is finite, and that $\pi^* \in \Pi$ achieves it:*

$$I_\ell(\pi^*; X) = \sup_{\pi \in \Pi} I_\ell(\pi; X) = C_\ell(\Pi) < \infty.$$

Let the divergence $D_\Omega(P_\theta, Q)$ be directionally lower semicontinuous in Q for each P_θ . Then the marginal distribution P_{π^} satisfies*

$$\sup_{\pi \in \Pi} L(P_{\pi^*}, \pi) \leq L(P_{\pi^*}, \pi^*) \leq \inf_Q L(Q, \pi^*),$$

so (P_{π^}, π^*) is a saddle point for the expected regret (19.3.5), and $L(P_{\pi^*}, \pi^*) = C_\ell(\Pi) = I_\ell(\pi^*; X)$. If additionally ℓ is strictly proper, then P_{π^*} uniquely achieves the infimum in $L(Q, \pi^*)$.*

Writing the lemma in terms of the gaps in the expected losses, if π^* maximizes $I_\ell(\pi; X)$ over $\pi \in \Pi$, then $Q^* = P_{\pi^*}$, and the regret/capacity game (19.0.3) has saddle point (P_{π^*}, π^*) :

$$\begin{aligned} \sup_{\pi \in \Pi} \int (\mathbb{E}_{P_\theta}[\ell(P_{\pi^*}, X)] - \mathbb{E}_{P_\theta}[\ell(P_\theta, X)]) d\pi(\theta) &\leq \int (\mathbb{E}_{P_\theta}[\ell(P_{\pi^*}, X)] - \mathbb{E}_{P_\theta}[\ell(P_\theta, X)]) d\pi^*(\theta) \\ &\leq \inf_Q \int (\mathbb{E}_{P_\theta}[\ell(Q, X)] - \mathbb{E}_{P_\theta}[\ell(P_\theta, X)]) d\pi^*(\theta). \end{aligned}$$

We temporarily defer the proof of Lemma 19.3.1 to Section 19.3.4.

Stating things informally for now, the following theorem captures what we view as prototypical regret/capacity duality:

Theorem 19.3.2 (Informal regret/capacity duality). *Let $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$ be a collection of distributions on \mathcal{X} , Π a convex collection of prior distributions on θ , and assume the loss ℓ is “proper enough” and that the capacity (19.3.4) is finite. Then there exists a unique distribution Q^* on \mathcal{X} such that*

$$\sup_{\pi \in \Pi} \int \text{EReg}(Q^*, P_\theta) d\pi(\theta) = C_\ell(\Pi).$$

In particular, the duality (19.3.7) holds, and

$$\inf_Q \sup_{\pi \in \Pi} \int \text{EReg}(Q, P_\theta) d\pi(\theta) = C_\ell(\Pi) = \sup_{\pi \in \Pi} \inf_Q \int \text{EReg}(Q, P_\theta) d\pi(\theta).$$

While we have been purposefully informal in the statement of Theorem 19.3.2, stating merely that the loss ℓ is “proper enough,” we provide several formalizations of the theorem in Section 19.6, presenting a full version in Theorem 19.6.5 in Section 19.6.4. In brief, “proper enough” means that the loss ℓ is strictly proper plus a bit more, so that convergence of distributions in the associated Bregman divergence representation of expected regret (19.3.1) guarantees some type of convergence of distributions.

In the next subsection, we give several instantiations of the theorem. From a practical perspective, Theorem 19.3.2 (and the more precise Theorem 19.6.5) could use some help: it only guarantees the existence of an optimal distribution achieving the capacity or worst-case expected regret. By returning to the setting in which we have sequential observations X_i , $i = 1, \dots, n$, we can ask not whether a worst-case regret minimizing distribution Q^* exists, but two related questions. First, whether we can achieve sublinear regret: there exists some distribution Q^* for which

$$\sup_{P \in \mathcal{P}} \text{EReg}_n(Q^*, P) \ll n.$$

Then the worst-case expected regret (or redundancy) grows only sublinearly with n , and we play (asymptotically) essentially as well as if we knew the true distribution P . Secondly and relatedly, we ask whether we can identify an (asymptotically) worst-case prior π^* . Because of loss propriety, we know that given a prior π , $Q = P_\pi$ minimizes $\int \text{EReg}(Q, P_\theta) d\pi(\theta)$, so that if we can perform estimation using the prior π^* , we achieve (asymptotically) optimal expected regret. These two questions form the core of Sections 19.4 and 19.5.

19.3.2 Instantiations and corollaries of regret/capacity duality

In spite of the informality of Theorem 19.3.2, we can still state rigorous consequences giving exemplar instantiations of the result. We do so here.

In the case of the log-loss, the regret is the redundancy (19.2.3), so that $\text{EReg}(Q, P) = D_{\text{kl}}(P\|Q)$. Then the information I_ℓ is simply the mutual information, and we have the classical capacity

$$C(\Pi) := \sup_{\pi} \{I(X; T) \mid \pi \text{ on } T \in \Theta, \pi \in \Pi\},$$

where the supremum is over a convex collection Π of distributions on T , and conditional on $T = \theta$ we draw $X \sim P_\theta$ as usual. We have the following *redundancy capacity* duality result, one of the foundational results in information theory.

Corollary 19.3.3 (Redundancy/capacity duality). *Let $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$ be distributions on a set \mathcal{X} , and assume that the capacity $C(\Pi) < \infty$. Then there exists a distribution Q^* on \mathcal{X} such that*

$$\sup_{\pi \in \Pi} \int D_{\text{kl}}(P_\theta\|Q^*) d\pi(\theta) = C(\Pi).$$

If $\Pi = \Pi(\Theta)$ consists of all distributions on Θ , then

$$\sup_{\theta \in \Theta} D_{\text{kl}}(P_\theta\|Q^*) = C(\Pi(\Theta)).$$

Lastly, if π^ achieves the supremum in the capacity, then $Q^* = P_{\pi^*}$ is the marginal over X when $T \sim \pi^*$, and we have the saddle point*

$$\sup_{\pi \in \Pi} \int D_{\text{kl}}(P_\theta\|P_{\pi^*}) d\pi(\theta) \leq \int D_{\text{kl}}(P_\theta\|P_{\pi^*}) d\pi^*(\theta) = I(X; T) \leq \inf_Q \int D_{\text{kl}}(P_\theta\|Q) d\pi^*(\theta).$$

Whenever \mathcal{X} is finite, we more or less have duality without conditions beyond ℓ being strictly proper and lower semicontinuous, that is, that $Q \mapsto \ell(Q, x)$ is lower semicontinuous.

Corollary 19.3.4 (Regret/capacity duality for finite sets). *Let $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$ be distributions on a finite set \mathcal{X} , let $\Pi \subset \Pi(\Theta)$ be a convex collection of distributions on Θ , and assume the capacity (19.3.4) is finite: $C_\ell(\Pi) < \infty$. If ℓ is strictly proper and lower semicontinuous, then there exists unique distribution Q^* on \mathcal{X} such that*

$$\sup_{\pi \in \Pi} \int (\mathbb{E}_{P_\theta}[\ell(Q^*, X)] - \mathbb{E}_{P_\theta}[\ell(P_\theta, X)]) d\pi(\theta) = C_\ell(\Pi).$$

If additionally $\Pi = \Pi(\Theta)$ consists of all distributions on Θ , then

$$\sup_{\theta \in \Theta} \mathbb{E}_{P_\theta}[\ell(Q^*, X)] - \mathbb{E}_{P_\theta}[\ell(P_\theta, X)] = C_\ell(\Pi(\Theta)).$$

Lastly, if π^ achieves the supremum in the capacity, then $Q^* = P_{\pi^*}$, and (P_{π^*}, π^*) is a saddle point.*

Finally, we state one more corollary of Theorem 19.6.5 to come, which shows an application of the results to the continuous-ranked probability scores we discuss in Examples 14.2.6 and 14.2.10, which arise when we predict cumulative distributions. Recall that in this case, for a CDF F , we have loss $\ell_{\text{crps}}(F, y) = \int (F(t) - \mathbf{1}\{y \leq t\})^2 dt$.

Corollary 19.3.5 (Regret/capacity duality for ranked scores). *Let $\mathcal{P} = \{F_\theta\}_{\theta \in \Theta}$ be any collection of cumulative distributions on a compact interval $T = [t_0, t_1]$ and let Π be a convex collection of distributions on Θ . Then there exists a unique cumulative distribution F^* such that*

$$\sup_{\pi \in \Pi} \int (\mathbb{E}_{F_\theta}[\ell_{\text{crps}}(F^*, X)] - \mathbb{E}_{F_\theta}[\ell(F_\theta, X)]) d\pi(\theta) = C_\ell(\Pi).$$

If π^ achieves the capacity $C_\ell(\Pi)$, then $F^*(t) = \int F_\theta(t) d\pi^*(t)$ is the marginal CDF.*

19.3.3 Maximum generalized entropy and Robust Bayesian procedures

The instantiations in Corollaries 19.3.3, 19.3.4, and 19.3.5 show the duality between expected regret (or redundancy) and capacity. We can specialize these results to show that there are robust Bayesian procedures for many of the loss-minimization games we originally discuss in Section 19.1. To demonstrate these dualities, we again modify our notation slightly, and consider the identity channel $X \rightarrow X$, that is, $T = X$. Then we assume that the losses are such that

$$\inf_P \ell(P, x) = 0 \text{ for all } x \in \mathcal{X}, \text{ i.e. } \ell(\mathbf{1}_x, x) = 0 \quad (19.3.9)$$

where $\mathbf{1}_x$ denotes the point mass on x . In this case, the (expected) regret is simply

$$\mathbb{E}_P[\ell(Q, X)] \geq 0,$$

and the worst case game is to solve

$$\text{minimize}_Q \sup_{P \in \mathcal{P}} \mathbb{E}_P[\ell(Q, X)]. \quad (19.3.10)$$

Example 19.3.6: If \mathcal{X} is discrete, then the log-loss $\ell(q, x) = -\log q(x)$ satisfies the condition (19.3.9), because $q = \mathbf{1}_x$ satisfies $\ell(q, x) = -\log 1 = 0$. \diamond

Example 19.3.7: If $\mathcal{X} \subset \mathbb{R}$, the continuous ranked probability score (CRPS) $\ell_{\text{crps}}(F, x) = \int (F(t) - \mathbf{1}\{x \leq t\})^2 dt$ satisfies the condition (19.3.9), where we take $F(t) = \mathbf{1}\{x \leq t\}$. \diamond

Assume now that \mathcal{P} is a convex collection of distributions on \mathcal{X} . Then recalling the generalized entropy (14.1.6)

$$H_\ell(P) = \inf_Q \mathbb{E}_P[\ell(Q, X)] = \mathbb{E}_P[\ell(P, X)],$$

we say that the robust game (19.3.10) has a solution if there exists Q^* such that

$$\sup_{P \in \mathcal{P}} H_\ell(P) = \sup_{P \in \mathcal{P}} \mathbb{E}_P[\ell(Q^*, X)],$$

and a saddle point (Q^*, P^*) if

$$\sup_{P \in \mathcal{P}} \mathbb{E}_P[\ell(Q^*, X)] \leq \mathbb{E}_{P^*}[\ell(Q^*, X)] \leq \inf_Q \mathbb{E}_{P^*}[\ell(Q, X)].$$

Note that if a saddle point exists, then whenever ℓ is strictly proper, it must be the case that $P^* = Q^*$, and moreover, P^* must maximize the generalized entropy $H_\ell(P)$ over \mathcal{P} .

We can therefore provide corollaries of Theorem 19.3.2 (or, more accurately, Theorem 19.6.5 to come).

Corollary 19.3.8. *Let \mathcal{X} be finite or countable, let the loss $\ell(q, x) = -\log q(x)$ be the log-loss, and let \mathcal{P} be a convex collection of distributions. For entropy $H(P) = -\sum_x p(x) \log p(x)$, assume that the maximum entropy $\sup_{P \in \mathcal{P}} H(P) < \infty$. Then*

$$\sup_{P \in \mathcal{P}} H(P) = \inf_Q \sup_{P \in \mathcal{P}} \mathbb{E}_P[-\log q(X)] < \infty,$$

and there exists a unique minimizer Q^ of $\sup_{P \in \mathcal{P}} \mathbb{E}_P[\ell(Q, X)]$. If P^* attains the maximum entropy, then $Q^* = P^*$ and (P^*, P^*) is a saddle point for the game (19.3.10).*

So maximum entropy distributions are robust Bayes.

To avoid annoying technicalities, we do not state the most general version of the results here. We can, however, give a few additional results in special cases, leaving generalizations to the exercises. Let us consider the case that space \mathcal{X} is finite. To give a more explicitly Bayesian flavor to the results, assume that we parameterize distributions on \mathcal{X} by some $\theta \in \Theta$ (this is no loss of generality), so that for any collection $\{P_\theta\}_{\theta \in \Theta}$, we have

$$\text{Conv}\{P_\theta\}_{\theta \in \Theta} = \{P_\pi\}_{\pi \in \Pi(\Theta)},$$

where $\Pi(\Theta)$ denotes the collection of prior probability distributions on Θ and P_π is the marginal over $T \sim \pi$ and $X \sim P_\theta$ conditional on $T = \theta$, as usual. (See Corollary B.1.16 in Appendix B.1.2 to see the equality of the convex hull.) Then for a given prior π , propriety of ℓ means that P_π minimizes $\int \mathbb{E}_{P_\theta}[\ell(Q, X)] d\pi(\theta) = \mathbb{E}_{P_\pi}[\ell(Q, X)]$. We call π^* a worst-case prior if it maximizes

$$\inf_Q \int \mathbb{E}_{P_\theta}[\ell(Q, X)] d\pi(\theta).$$

Corollary 19.3.9. *Let \mathcal{X} be finite and ℓ be strictly proper. Let $\{P_\theta\}_{\theta \in \Theta}$ be any collection of distributions on \mathcal{X} . Assume the maximum generalized entropy $\sup_{\pi \in \Pi(\Theta)} H_\ell(P_\pi) < \infty$. Then there exists a unique minimizer Q^* of $\sup_{P \in \mathcal{P}} \mathbb{E}_P[\ell(Q, X)]$, and*

$$\sup_{\pi \in \Pi(\Theta)} H_\ell(P_\pi) = \sup_{\pi \in \Pi(\Theta)} \mathbb{E}_{P_\pi}[\ell(Q^*, X)] = \inf_Q \sup_{\pi \in \Pi(\Theta)} \int \mathbb{E}_{P_\theta}[\ell(Q, X)] d\pi(\theta) < \infty.$$

If additionally the convex hull $\{P_\pi\}_{\pi \in \Pi(\Theta)}$ is closed, then there is a prior π^ maximizing the generalized entropy $H_\ell(P_\pi)$, π^* is the worst-case prior, $Q^* = P_{\pi^*}$, and (P_{π^*}, P_{π^*}) is a saddle point for the game (19.3.10).*

Proof All that we need to argue is that there is a prior π^* maximizing the entropy $H_\ell(P_\pi)$. Because $\mathcal{P} = \text{Conv}\{P_\theta\}_{\theta \in \Theta} \subset \Delta_{\mathcal{X}}$ is bounded, if it is closed it is necessarily compact. Then noting that $H_\ell(P)$ is the infimum of linear functions of P , it is concave and upper semi-continuous, so that its maximizers are attained over compact sets. So there is $P^* \in \mathcal{P}$ maximizing $H_\ell(P)$; that \mathcal{P} is closed implies there exists at least one prior π^* on Θ for which $P^* = P_{\pi^*}$. \square

We have no particular need to work with finite domains, though (as we note above), this usually necessitates more care in the constructions. As a particular example, however, we can leverage Corollary 19.3.5.

Corollary 19.3.10. *Let $T = [t_0, t_1] \subset \mathbb{R}$ be a compact interval, and let \mathcal{P} be the collection of cumulative distributions on T (i.e., distributions P supported on T). Then for the continuous ranked probability loss ℓ_{crps} , there is a (unique) CDF F^* minimizing $\sup_{F \in \mathcal{P}} \mathbb{E}_F[\ell_{\text{crps}}(F, X)]$, and*

$$\sup_{F \in \mathcal{P}} \int_{t_0}^{t_1} (1 - F(t))F(t)dt = \sup_{F \in \mathcal{P}} \mathbb{E}_F[\ell_{\text{crps}}(F^*, X)].$$

JCD Comment: Should probably add some kind of exercises on these

19.3.4 Proof of Lemma 19.3.1

Because ℓ is proper, for any π , the distribution $P_\pi := \int P_\theta d\pi(\theta)$ minimizes $L(Q, \pi)$ over distributions Q (recall Eq. (19.3.6)). Let $Q^* = P_{\pi^*}$. Then for any distribution Q , we have

$$C_\ell = I_\ell(\pi^*; X) = \inf_Q L(Q, \pi^*) \leq L(Q, \pi^*). \quad (19.3.11)$$

Now, take any distribution $\pi \in \Pi$, and let $\lambda \in (0, 1)$, defining $\pi_\lambda = (1 - \lambda)\pi^* + \lambda\pi$ (recall we have assumed Π is convex). Then $I_\ell(\pi_\lambda; X) \leq C_\ell$ immediately, while for $P_\lambda = \lambda P_\pi + (1 - \lambda)P_{\pi^*}$ and $T \in \Theta$ representing a draw from π or π^* , we have

$$\begin{aligned} C_\ell &\geq I_\ell(\pi_\lambda; X) = \inf_Q \left\{ \int \mathbb{E}_{P_\theta}[\ell(Q, X)] d\pi_\lambda(\theta) - \lambda \mathbb{E}_\pi[H_\ell(P_T)] - (1 - \lambda) \mathbb{E}_{\pi^*}[H_\ell(P_T)] \right\} \\ &\stackrel{(i)}{=} \int \mathbb{E}_{P_\theta}[\ell(P_\lambda, X)] d\pi_\lambda(\theta) - \lambda \mathbb{E}_\pi[H_\ell(P_T)] - (1 - \lambda) \mathbb{E}_{\pi^*}[H_\ell(P_T)] \\ &= (1 - \lambda) \left[\int \mathbb{E}_{P_\theta}[\ell(P_\lambda, X)] d\pi^*(\theta) - \mathbb{E}_{\pi^*}[H_\ell(P_T)] \right] + \lambda \left[\int \mathbb{E}_{P_\theta}[\ell(P_\lambda, X)] d\pi(\theta) - \mathbb{E}_\pi[H_\ell(P_T)] \right] \\ &\stackrel{(ii)}{\geq} (1 - \lambda)C + \lambda \left(\int \mathbb{E}_{P_\theta}[\ell(P_\lambda, X)] d\pi(\theta) - \mathbb{E}_\pi[H_\ell(P_T)] \right), \end{aligned}$$

where step (i) uses the propriety of the loss and step (ii) follows from the earlier bound (19.3.11). Rearranging and dividing through by $\lambda > 0$, we obtain by definition of L that

$$L(\lambda P_\pi + (1 - \lambda)P_{\pi^*}, \pi) = \int \mathbb{E}_{P_\theta}[\ell(\lambda P_\pi + (1 - \lambda)P_{\pi^*}, X)] d\pi(\theta) - \mathbb{E}_\pi[H_\ell(P_T)] \leq C_\ell$$

for all $\pi \in \Pi$ and all $\lambda \in (0, 1)$.

Finally, we use that the expected regret $D_\Omega(P_\theta, Q) = \mathbb{E}_{P_\theta}[\ell(Q, X)] - \mathbb{E}_{P_\theta}[\ell(P_\theta, X)]$ is lower semicontinuous along one-dimensional versions of Q by assumption. This in particular implies that

$$\liminf_{\lambda \downarrow 0} \mathbb{E}_{P_\theta}[\ell(\lambda P + (1 - \lambda)Q, X)] - \mathbb{E}_{P_\theta}[\ell(P_\theta, X)] \geq D_\Omega(P_\theta, Q)$$

for any P, Q . Applying Fatou's lemma (see Proposition A.3.1 in Appendix A.3) then gives

$$\begin{aligned} L(P_{\pi^*}, \pi) &\leq \int \liminf_{\lambda \downarrow 0} (\mathbb{E}_{P_\theta}[\ell(\lambda P_\pi + (1 - \lambda)P_{\pi^*}, X)] - H_\ell(P_\theta)) d\pi(\theta) \\ &\leq \liminf_{\lambda \downarrow 0} L(\lambda P_\pi + (1 - \lambda)P_{\pi^*}, \pi) \leq C_\ell. \end{aligned}$$

This implies the desired saddle point.

19.4 Minimax strategies for regret

With general theory about the existence of saddle points and minimax solutions in place, we now turn to more explicit strategies for attaining optimal regret as the sample size $n \rightarrow \infty$. We discuss the adversarial setting only briefly, as optimal strategies are somewhat difficult to implement; the redundancy setting allows easier exploration. We focus on the log loss, because explicit descriptions of optimal and near-optimal strategies for other losses often remain open questions.

We begin by describing a notion of complexity that captures the best possible regret in the adversarial setting. Assume without loss of generality that we have a set of distributions $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$ parameterized by $\theta \in \Theta$, where the distributions are supported on \mathcal{X}^n . We define the complexity of the set \mathcal{P} (viz. the complexity of Θ) as

$$\text{Comp}_n(\Theta) := \log \int_{\mathcal{X}^n} \sup_{\theta \in \Theta} p_\theta(x_1^n) dx_1^n \quad \text{or} \quad \text{Comp}_n(\Theta) := \log \int_{\mathcal{X}^n} \sup_{\theta \in \Theta} p_\theta(x_1^n) d\nu(x_1^n), \quad (19.4.1)$$

where ν is some base measure on \mathcal{X}^n . Note that we may have $\text{Comp}_n(\Theta) = +\infty$, especially when Θ is non-compact. This is not particularly uncommon, for example, consider the case of a normal location family model over $\mathcal{X} = \mathbb{R}$ with $\Theta = \mathbb{R}$.

The complexity (19.4.1) is precisely the minimax regret in the adversarial setting.

Proposition 19.4.1. *The minimax regret*

$$\inf_Q \mathfrak{R}^{\mathcal{X}}(Q, \mathcal{P}) = \text{Comp}_n(\Theta).$$

Moreover, if $\text{Comp}_n(\Theta) < +\infty$, then the normalized maximum likelihood distribution (also known as the Shtarkov distribution) \bar{Q} , defined with density

$$\bar{q}(x_1^n) = \frac{\sup_{\theta \in \Theta} p_\theta(x_1^n)}{\int \sup_{\theta} p_\theta(x_1^n) dx_1^n},$$

is uniquely minimax optimal.

The proposition completely characterizes the minimax regret in the adversarial setting, and it gives the unique distribution achieving the regret. Unfortunately, in most cases it is challenging to compute the minimax optimal distribution \bar{Q} , so we must make approximations of some type. In the sequel, we will see that by moving to redundancy rather than adversarial regret, Bayesian approximations to \bar{Q} allow this.

Proof We begin by proving the result in the case that $\text{Comp}_n < +\infty$. First, note that the normalized maximum likelihood distribution \bar{Q} has constant regret:

$$\begin{aligned} \mathfrak{R}_n^{\mathcal{X}}(\bar{Q}, \mathcal{P}) &= \sup_{x_1^n \in \mathcal{X}^n} \left[\log \frac{1}{\bar{q}(x_1^n)} - \log \frac{1}{\sup_{\theta} p_\theta(x_1^n)} \right] \\ &= \sup_{x_1^n} \left[\log \frac{\int \sup_{\theta} p_\theta(x_1^n) dx_1^n}{\sup_{\theta} p_\theta(x_1^n)} - \log \frac{1}{\sup_{\theta} p_\theta(x_1^n)} \right] = \text{Comp}_n(\mathcal{P}). \end{aligned}$$

Moreover, for any distribution Q on \mathcal{X}^n we have

$$\begin{aligned} \mathfrak{R}_n^{\mathcal{X}}(Q, \mathcal{P}) &\geq \int \left[\log \frac{1}{q(x_1^n)} - \log \frac{1}{\sup_{\theta} p_\theta(x_1^n)} \right] \bar{q}(x_1^n) dx_1^n \\ &= \int \left[\log \frac{\bar{q}(x_1^n)}{q(x_1^n)} + \text{Comp}_n(\Theta) \right] \bar{q}(x_1^n) dx_1^n \\ &= D_{\text{kl}}(\bar{Q} \| Q) + \text{Comp}_n(\Theta), \end{aligned} \quad (19.4.2)$$

so that \bar{Q} is uniquely minimax optimal, as $D_{\text{kl}}(\bar{Q}\|Q) > 0$ unless $\bar{Q} = Q$.

Now we show how to extend the lower bound (19.4.2) to the case when $\text{Comp}_n(\Theta) = +\infty$. Let us assume without loss of generality that \mathcal{X} is countable and consists of points x_1, x_2, \dots (we can discretize \mathcal{X} otherwise) and assume we have $n = 1$. Fix any $\epsilon \in (0, 1)$ and construct the sequence $\theta_1, \theta_2, \dots$ so that $p_{\theta_j}(x_j) \geq (1 - \epsilon) \sup_{\theta \in \Theta} p_{\theta}(x)$, and define the sets $\Theta_j = \{\theta_1, \dots, \theta_j\}$. Clearly we have $\text{Comp}(\Theta_j) \leq \log j$, and if we define $\bar{q}_j(x) = \max_{\theta \in \Theta_j} p_{\theta}(x) / \sum_{x \in \mathcal{X}} \max_{\theta \in \Theta_j} p_{\theta}(x)$, we may extend the reasoning yielding inequality (19.4.2) to obtain

$$\begin{aligned} \mathfrak{R}^{\mathcal{X}}(Q, \mathcal{P}) &= \sup_{x \in \mathcal{X}} \left[\log \frac{1}{q(x)} - \log \frac{1}{\sup_{\theta \in \Theta} p_{\theta}(x)} \right] \\ &\geq \sum_x \bar{q}_j(x) \left[\log \frac{1}{q(x)} - \log \frac{1}{\max_{\theta \in \Theta_j} p_{\theta}(x)} \right] \\ &= \sum_x \bar{q}_j(x) \left[\log \frac{\bar{q}_j(x)}{q(x)} + \log \sum_{x'} \max_{\theta \in \Theta_j} p_{\theta}(x') \right] = D_{\text{kl}}(\bar{Q}_j\|Q) + \text{Comp}(\Theta_j). \end{aligned}$$

But of course, by noting that

$$\text{Comp}(\Theta_j) \geq (1 - \epsilon) \sum_{i=1}^j \sup_{\theta} p_{\theta}(x_i) + \sum_{i>j} \max_{\theta \in \Theta_j} p_{\theta}(x_i) \rightarrow +\infty$$

as $j \rightarrow \infty$, we obtain the result when $\text{Comp}_n(\Theta) = \infty$. \square

While typically computing the minimax (adversarial) regret is challenging, in some simple cases we can compute it to within constant factors. In this case, we compete with the family of i.i.d. Bernoulli distributions.

Example 19.4.2 (Complexity of the Bernoulli distribution): Consider competing against the family of Bernoulli distributions $\{P_{\theta}\}_{\theta \in [0,1]}$, where for a point $x \in \{0, 1\}$, we have $P_{\theta}(x) = \theta^x(1 - \theta)^{1-x}$. For a sequence $x_1^n \in \{0, 1\}^n$ with m non-zeros, we thus have for $\hat{\theta} = m/n$ that

$$\sup_{\theta \in [0,1]} P_{\theta}(x_1^n) = P_{\hat{\theta}}(x_1^n) = \hat{\theta}^m(1 - \hat{\theta})^{n-m} = \exp(-nh_2(\hat{\theta})),$$

where $h_2(p) = -p \log p - (1 - p) \log(1 - p)$ is the binary entropy. Using this representation, we find that the complexity of the Bernoulli family is

$$\text{Comp}_n([0, 1]) = \log \sum_{m=0}^n \binom{n}{m} e^{-nh_2(\frac{m}{n})}.$$

Rather than explicitly compute with this, we now use Stirling's approximation (cf. Cover and Thomas [57, Lemma 17.5.1]): for any $p \in (0, 1)$ with $np \in \mathbb{N}$, we have

$$\binom{n}{np} \in \frac{1}{\sqrt{n}} \left[\frac{1}{\sqrt{8p(1-p)}}, \frac{1}{\sqrt{\pi p(1-p)}} \right] \exp(nh_2(p)).$$

Thus, by dealing with the boundary cases $m = n$ and $m = 0$ explicitly, we obtain

$$\begin{aligned} \sum_{m=0}^n \binom{n}{m} \exp\left(-nh_2\left(\frac{m}{n}\right)\right) &= 2 + \sum_{m=1}^{n-1} \binom{n}{m} \exp\left(-nh_2\left(\frac{m}{n}\right)\right) \\ &\in 2 + \left[\frac{1}{\sqrt{8}}, \frac{1}{\sqrt{\pi}}\right] \frac{1}{\sqrt{n}} \underbrace{\sum_{m=1}^{n-1} \frac{1}{\sqrt{\frac{m}{n}(1-\frac{m}{n})}}}_{\rightarrow n \int_0^1 (\theta(1-\theta))^{-\frac{1}{2}}}, \end{aligned}$$

the noted asymptote occurring as $n \rightarrow \infty$ because this sum is a Riemann sum for the integral $\int_0^1 \theta^{-1/2}(1-\theta)^{-1/2} d\theta$. So as $n \rightarrow \infty$,

$$\begin{aligned} \inf_Q \mathfrak{R}_n^{\mathcal{X}}(Q, \mathcal{P}) &= \text{Comp}_n([0, 1]) = \log \left(2 + [8^{-1/2}, \pi^{-1/2}] n^{1/2} \int_0^1 \frac{1}{\sqrt{\theta(1-\theta)}} d\theta \right) + o(1) \\ &= \frac{1}{2} \log n + \log \int_0^1 \frac{1}{\sqrt{\theta(1-\theta)}} d\theta + O(1). \end{aligned}$$

We remark in passing that this is equal to $\frac{1}{2} \log n + \log \int_0^1 \sqrt{J_\theta} d\theta$, where J_θ denotes the Fisher information of the Bernoulli family (recall Example 3.1.1). We will see that this holds in more generality, at least for redundancy, in the sequel. \diamond

19.5 Mixture (Bayesian) strategies and redundancy

In the less adversarial setting of expected regret, we compete against a random sequence X_1^n of data, drawn from some fixed distribution P , rather than an adversarially chosen sequence x_1^n . Thinking of this problem as a game, we choose a distribution Q according to which we make predictions (based on previous data), and nature chooses a distribution $P_\theta \in \mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$. Simplifying the setting (19.2.3) to make things tractable, we let the data X_1^n be generated i.i.d. according to P_θ , and we suffer expected regret (or redundancy)

$$\text{Red}_n(Q, P_\theta) = \mathbb{E}_\theta \left[\log \frac{1}{q(X_1^n)} \right] - \mathbb{E}_\theta \left[\log \frac{1}{p_\theta(X_1^n)} \right] = D_{\text{kl}}(P_\theta^n \| Q_n), \quad (19.5.1)$$

where we use Q_n to denote that Q is applied on all n data points (in a sequential fashion, as $Q(\cdot | X_1^{i-1})$). In this expression, q and p denote the densities of Q and P , respectively. In a slightly more general setting, we may consider the expected regret of Q with respect to a distribution P_θ even under model mis-specification, meaning that the data is generated according to an alternate distribution P . In this case, the (more general) redundancy becomes

$$\mathbb{E}_P \left[\log \frac{1}{q(X_1^n)} - \log \frac{1}{p_\theta(X_1^n)} \right]. \quad (19.5.2)$$

In both cases (19.5.1) and (19.5.2), we would like to be able to guarantee that the redundancy grows more slowly than n as $n \rightarrow \infty$. That is, we would like to find distributions Q such that, for any $\theta_0 \in \Theta$, we have $\frac{1}{n} D_{\text{kl}}(P_{\theta_0}^n \| Q_n) \rightarrow 0$ as $n \rightarrow \infty$. Assuming we could actually obtain such a distribution in general, this is interesting because (even in the i.i.d. case) for *any* fixed distribution $P_\theta \neq P_{\theta_0}$, we must have $D_{\text{kl}}(P_{\theta_0}^n \| P_\theta^n) = n D_{\text{kl}}(P_{\theta_0} \| P_\theta) = \Omega(n)$. Motivated by the

redundancy/capacity duality Corollary 19.3.3, a natural approach is to consider mixture distributions, where we choose Q as a convex combination (mixture) of all the possible source distributions P_θ for $\theta \in \Theta$, that is, we let $Q = P_\pi^n$ for some prior π on Θ . Here, in contrast to the exact dualities in Section 19.3, we will not seek explicit worst-case priors, instead relying on asymptotics to show approximate optimality.

To see how we make predictions from such mixtures, let $Q_\pi^n(A) = \int \pi(\theta) P_\theta^n(A) d\theta$ for $A \subset \mathcal{X}^n$, where we note that Q_π^n will typically not be a product distribution. Then at time i , this mixture plays the density

$$q_\pi(x_i | x_1^{i-1}) = \int_{\Theta} q(x_i, \theta | x_1^{i-1}) d\theta = \int_{\Theta} p_\theta(x_i) \pi(\theta | x_1^{i-1}) d\theta$$

by construction of the distributions Q_π as mixtures of i.i.d. P_θ . Here, the posterior distribution $\pi(\theta | x_1^{i-1})$ takes the form

$$\pi(\theta | x_1^{i-1}) = \frac{\pi(\theta) p_\theta(x_1^{i-1})}{\int_{\Theta} \pi(\theta') p_{\theta'}(x_1^{i-1}) d\theta'} = \frac{\pi(\theta) \exp\left(-\log \frac{1}{p_\theta(x_1^{i-1})}\right)}{\int_{\Theta} \pi(\theta') p_{\theta'}(x_1^{i-1}) d\theta'}, \quad (19.5.3)$$

where we have emphasized that this strategy exhibits an *exponential weighting* approach, where distribution weights are scaled exponentially by their previous loss performance of $\log 1/p_\theta(x_1^{i-1})$.

The weighting scheme (19.5.3) asymptotically strong performance. In fact, we say that so long as the prior π puts non-zero mass over all of Θ , under some appropriate smoothness conditions, the scheme Q_π is universal, meaning that $D_{\text{kl}}(P_\theta^n \| Q_\pi^n) = o(n)$ for any $\theta \in \text{int } \Theta$. Clarke and Barron [52, 53] make these statements rigorous; while the particular conditions are beyond the scope of this book, in essence they require the following: the Fisher information J_θ for the family $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$ exists in a compact set K interior to Θ , and the distributions P_θ are sufficiently regular that differentiation and integration can be interchanged (essentially, uniform versions of the conditions (i)–(v) necessary for the van Trees inequality, Theorem 12.2.7).

Theorem 19.5.1 (Clarke and Barron [52, 53]). *Let $Q_\pi^n = \int P_\theta^n \pi(\theta) d\theta$ be the marginal distribution over X_1^n obtained by drawing $T \sim \pi$ and then $X_i \stackrel{\text{iid}}{\sim} P_\theta$ conditional on $T = \theta$. Let π have compact support Θ , and $K \subset \text{int } \Theta$ be any compact set. Then under appropriate smoothness conditions on the distributions P_θ and π ,*

$$D_{\text{kl}}(P_\theta^n \| Q_\pi^n) - \frac{d}{2} \log \frac{n}{2\pi e} \rightarrow \log \frac{1}{\pi(\theta)} + \frac{1}{2} \log \det(J_\theta) \quad \text{as } n \rightarrow \infty \quad (19.5.4a)$$

uniformly in $\theta \in K$, and

$$\left| \int D_{\text{kl}}(P_\theta^n \| Q_\pi^n) \pi(\theta) d\theta - \frac{d}{2} \log \frac{n}{2\pi e} - \int \pi(\theta) \log \frac{\sqrt{\det(J_\theta)}}{\pi(\theta)} d\theta \right| \rightarrow 0. \quad (19.5.4b)$$

While we do not rigorously prove the theorem, we give a sketch showing the main components of the result based on asymptotic normality arguments in Section 19.5.2.

Example 19.5.2 (Bernoulli distributions with a Beta prior): Consider the class of binary (i.i.d. or memoryless) Bernoulli sources, that is, the X_i are i.i.d Bernoulli(θ), where $\theta = P_\theta(X = 1) \in [0, 1]$. The Beta(α, β)-distribution prior on θ is the mixture π with density

$$\pi(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

on $[0, 1]$, where $\Gamma(a) = \int_0^\infty t^{a-1} e^{-t} dt$ denotes the gamma function. We remark that that under the **Beta** (α, β) distribution, we have $\mathbb{E}_\pi[\theta] = \frac{\alpha}{\alpha+\beta}$. (See any undergraduate probability text for such results.)

If we play via a mixture of Bernoulli distributions under such a **Beta**-prior for θ , by Theorem 19.5.1 we have a universal prediction scheme. We may also explicitly calculate the predictive distribution Q . To do so, we first compute the posterior $\pi(\theta \mid X_1^i)$ as in expression (19.5.3). Let $S_i = \sum_{j=1}^i X_j$ be partial sum of the X s up to iteration i . Then

$$\pi(\theta \mid x_1^i) = \frac{p_\theta(x_1^i)\pi(\theta)}{q(x_1^i)} \propto \theta^{S_i}(1-\theta)^{i-S_i}\theta^{\alpha-1}\theta^{\beta-1} = \theta^{\alpha+S_i-1}(1-\theta)^{\beta+i-S_i-1},$$

where we have ignored the denominator as we must simply normalize the above quantity in θ . But by inspection, the posterior density of $\theta \mid X_1^i$ is a **Beta** $(\alpha + S_i, \beta + i - S_i)$ distribution. Thus to compute the predictive distribution, we note that $\mathbb{E}_\theta[X_i] = \theta$, so we have

$$Q(X_i = 1 \mid X_1^i) = \mathbb{E}_\pi[\theta \mid X_1^i] = \frac{S_i + \alpha}{i + \alpha + \beta}.$$

Moreover, Theorem 19.5.1 shows that when we play the prediction game with a **Beta** (α, β) -prior, we have redundancy scaling as

$$D_{\text{kl}}(P_{\theta_0}^n \| Q_\pi^n) = \frac{1}{2} \log \frac{n}{2\pi e} + \log \left[\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \frac{1}{\theta_0^{\alpha-1}(1-\theta_0)^{\beta-1}} \right] + \frac{1}{2} \log \frac{1}{\theta_0(1-\theta_0)} + o(1)$$

for $\theta_0 \in (0, 1)$. \diamond

We can also show how to play sequentially with Gaussian models, a simplified version of Example 19.2.3.

Example 19.5.3 (Gaussian distributions with a Gaussian prior): Consider the collection of i.i.d. (memoryless) Gaussian sources, where $X_i \stackrel{\text{iid}}{\sim} \mathbf{N}(\theta, \sigma^2 I)$. Then if the prior π on θ is $\mathbf{N}(0, \tau^2 I)$, a calculation involving completing squares yields posterior density

$$\pi(\theta \mid x_1^n) \propto \exp \left(-\frac{1}{2} \left(\frac{1}{\tau^2} + \frac{n}{\sigma^2} \right) \left\| \theta - \frac{n\tau^2}{n\tau^2 + \sigma^2} \bar{x}_n \right\|_2^2 \right),$$

that is, if $X_i \mid T = \theta \sim \mathbf{N}(\theta, \sigma^2 I)$ and $T \sim \mathbf{N}(0, \tau^2 I)$,

$$T \mid X_1^n \sim \mathbf{N} \left(\frac{\tau^2 n}{\tau^2 n + \sigma^2} \bar{X}_n, \frac{\sigma^2 \tau^2}{\sigma^2 + \tau^2 n} I \right).$$

Thus, the predictive distribution $Q_\pi(X_{n+1} \in \cdot \mid X_1^n)$ is Gaussian with mean $\mathbb{E}_Q[X_{n+1} \mid X_1^n] = \mathbb{E}_\pi[T \mid X_1^n] = \frac{\tau^2 n}{\tau^2 n + \sigma^2} \bar{X}_n$ and covariance

$$\text{Cov}(X_{n+1} \mid X_1^n) = \text{Cov}(T \mid X_1^n) + \sigma^2 I = \left(\frac{\sigma^2 \tau^2}{\sigma^2 + \tau^2 n} + \sigma^2 \right) I.$$

Direct, but tedious, calculations show that this agrees with the limit (19.5.4a). \diamond

While Theorem 19.5.1 addresses the case in which the data are generated i.i.d. P_θ , so that the model P_θ is well-specified, mixture models enjoy a type of robustness even under model misspecification, that is, when the true distribution generating the data does not belong to the class $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$. In this case, we look at the generalized redundancy (19.5.2), measuring loss relative to $D_{\text{kl}}(P \| P_\theta)$. We now allow restricting mixture distributions to a subset $\Theta_0 \subset \Theta$ by defining

$$Q_{\pi, \Theta_0}(A) := \frac{1}{\pi(\Theta_0)} \int_{\Theta_0} P_\theta(A) d\pi(\theta).$$

Then we obtain the following robustness result.

Proposition 19.5.4. *Assume that P_θ have densities p_θ over \mathcal{X} , let P be any distribution having density p over \mathcal{X} , and let q_π be the density associated with Q_π . Then for all $\Theta_0 \subset \Theta$,*

$$\mathbb{E}_P \left[\log \frac{1}{q_\pi(X)} - \log \frac{1}{p_\theta(X)} \right] \leq \log \frac{1}{\pi(\Theta_0)} + D_{\text{kl}}(P \| Q_{\pi, \Theta_0}) - D_{\text{kl}}(P \| P_\theta).$$

In particular, Proposition 19.5.4 shows that so long as the mixture distributions Q_{π, Θ_0} can closely approximate P_θ , then we attain a convergence guarantee nearly as good as any in the family $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$. (This result is similar in flavor to the mutual information bound (10.2.2), Corollary 10.2.2, and the *index of resolvability* there.)

Proof Fix any $\Theta_0 \subset \Theta$. Then $q_\pi(x) = \int_{\Theta} p_\theta(x) d\pi(\theta) \geq \int_{\Theta_0} p_\theta(x) d\pi(\theta)$. Thus we have

$$\begin{aligned} \mathbb{E}_P \left[\log \frac{p(X)}{q_\pi(X)} \right] &\leq \mathbb{E}_P \left[\inf_{\Theta_0 \subset \Theta} \log \frac{p(X)}{\int_{\Theta_0} p_\theta(x) d\pi(\theta)} \right] \\ &= \mathbb{E}_P \left[\inf_{\Theta_0} \log \frac{p(X) \pi(\Theta_0)}{\pi(\Theta_0) \int_{\Theta_0} p_\theta(x) d\pi(\theta)} \right] = \mathbb{E}_P \left[\inf_{\Theta_0} \log \frac{p(X)}{\pi(\Theta_0) q_{\pi, \Theta_0}(X)} \right]. \end{aligned}$$

This is certainly smaller than the same quantity with the infimum outside the expectation, and noting that

$$\mathbb{E}_P \left[\log \frac{1}{q_\pi(X)} - \log \frac{1}{p_\theta(X)} \right] = \mathbb{E}_P \left[\log \frac{p(X)}{q_\pi(X)} \right] - \mathbb{E}_P \left[\log \frac{p(X)}{p_\theta(X)} \right]$$

gives the result. \square

19.5.1 Bayesian redundancy and objective, reference, and Jeffreys priors

Consider again the Bayesian redundancy

$$\inf_Q \int_{\Theta} \pi(\theta) D_{\text{kl}}(P_\theta^n \| Q) d\theta = \int_{\Theta} \pi(\theta) D_{\text{kl}}(P_\theta \| Q_\pi^n) d\theta = I_\pi(T; X_1^n),$$

the mutual information between a random variable $T \sim \pi$ and $X_1^n \stackrel{\text{iid}}{\sim} P_\theta$ conditional on $T = \theta$. With Theorem 19.5.1 in hand, we can give a somewhat more nuanced picture of this mutual information quantity. As a first consequence of Theorem 19.5.1, we have that

$$I_\pi(T; X_1^n) = \frac{d}{2} \log \frac{n}{2\pi e} + \int \pi(\theta) \log \frac{\sqrt{\det J_\theta}}{\pi(\theta)} d\theta + o(1), \quad (19.5.5)$$

where J_θ denotes the Fisher information matrix for the family $\{P_\theta\}_{\theta \in \Theta}$. One strand of Bayesian statistics known as *reference analysis* advocates that in performing a Bayesian analysis, that is, performing an experiment (observing X_1^n) to estimate $\theta \in \Theta$, we should choose the prior π that maximizes the mutual information between the parameters θ about which we wish to make inferences and any observations X_1^n available. Moreover, in this set of strategies, one allows n to tend to ∞ , as we wish to take advantage of any data we might actually see. The asymptotic formula (19.5.5) allows us to choose such a prior. (We will not delve too deeply into this; instead, we refer to the survey by Bernardo [28] and papers of Berger et al. [25, 26].)

In a different vein, Jeffreys [121] proposed that if the square root of the determinant of the Fisher information was integrable, then one should take π as

$$\pi_{\text{jeffreys}}(\theta) = \frac{\sqrt{\det J_\theta}}{\int_{\Theta} \sqrt{\det J_\theta} d\theta}$$

giving the eponymous *Jeffreys prior*. Jeffreys originally proposed this for invariance reasons, as the inferences made on the parameter θ under the prior π_{jeffreys} are identical to those made on a transformed parameter $\phi(\theta)$ under the appropriately transformed Jeffreys prior.

Proceeding somewhat non-rigorously, in that we shall not worry about integrability, or uniformity in the prior π , the asymptotic expression (19.5.5) shows that the Jeffreys prior and the asymptotic reference prior coincide. Indeed, computing the integral in (19.5.5), we have

$$\begin{aligned} \int_{\Theta} \pi(\theta) \log \frac{\sqrt{\det J_\theta}}{\pi(\theta)} d\theta &= \int_{\Theta} \pi(\theta) \log \frac{\pi_{\text{jeffreys}}(\theta)}{\pi(\theta)} d\theta + \log \int \sqrt{\det J_\theta} d\theta \\ &= -D_{\text{kl}}(\pi \| \pi_{\text{jeffreys}}) + \log \int \sqrt{\det J_\theta} d\theta, \end{aligned}$$

whenever the Jeffreys prior exists. Moreover, we see that in an asymptotic sense, the Jeffreys prior is the worst-case prior distribution π for nature to play, as otherwise the $-D_{\text{kl}}(\pi \| \pi_{\text{jeffreys}})$ term in the expected (Bayesian) redundancy is negative.

Example 19.5.5 (Jeffreys priors and the exponential distribution): Let us now assume that our source distributions P_θ are exponential distributions, meaning that $\theta \in (0, \infty)$ and we have density $p_\theta(x) = \exp(-\theta x - \log \frac{1}{\theta})$ for $x \in [0, \infty)$. This is clearly an exponential family model, and the Fisher information is easy to compute as $J_\theta = \frac{\partial^2}{\partial \theta^2} \log \frac{1}{\theta} = 1/\theta^2$ (cf. Eq. (3.3.2)).

In this case, the Jeffreys prior is $\pi_{\text{jeffreys}}(\theta) \propto \sqrt{J} = 1/\theta$, but this “density” does not integrate over $[0, \infty)$. One approach to this difficulty, advocated by Bernardo [28, Definition 3] (among others) is to just proceed formally and notice that after observing a single datapoint, the “posterior” distribution $\pi(\theta | X)$ is well-defined. Following this idea, note that after seeing some data X_1, \dots, X_i , with $S_i = \sum_{j=1}^i X_j$ as the partial sum, we have

$$\pi(\theta | x_1^i) \propto p_\theta(x_1^i) \pi_{\text{jeffreys}}(\theta) = \theta^i \exp\left(-\theta \sum_{j=1}^i x_j\right) \frac{1}{\theta} = \theta^{i-1} \exp(-\theta S_i).$$

Integrating, we have for $s_i = \sum_{j=1}^i x_j$

$$q(x | x_1^i) = \int_0^\infty p_\theta(x) \pi(\theta | x_1^i) d\theta \propto \int_0^\infty \theta e^{-\theta x} \theta^{i-1} e^{-\theta s_i} d\theta = \frac{1}{(s_i + x)^{i+1}} \int_0^\infty u^i e^{-u} du,$$

where we made the change of variables $u = \theta(s_i + x)$. This is at least a distribution that normalizes, so often one simply assumes the existence of a piece of fake data. For example, by saying we “observe” $x_0 = 1$, we have prior proportional to $\pi(\theta) = e^{-\theta}$, which yields redundancy

$$D_{\text{kl}}(P_{\theta_0}^n \| Q_n^\pi) = \frac{1}{2} \log \frac{n}{2\pi e} + \theta_0 + \log \frac{1}{\theta_0} + o(1).$$

The difference is that, in this case, the redundancy bound is no longer uniform in θ_0 , as it would be for the true reference (or Jeffreys, if it exists) prior. \diamond

19.5.2 Heuristic calculations: normality and Theorem 19.5.1

In this section, we very briefly (and very hand-wavily) justify the asymptotic expression (19.5.4a), from which the expectation version (19.5.4b) more or less follows. To do this, we argue that the posterior distribution $\pi(\theta | X_1^n)$ should be approximately normally distributed with appropriate variance measure, which gives the result. (Clarke and Barron [53] provide a fully rigorous proof.)

Fix θ_0 . We begin by expanding the log likelihood in the divergence

$$D_{\text{kl}}(P_{\theta_0} \| Q_\pi^n) = \mathbb{E}_{P_{\theta_0}} \left[\log \frac{p_{\theta_0}(X_1^n)}{q_\pi(X_1^n)} \right] = \mathbb{E}_{P_{\theta_0}} \left[\log \frac{p_{\theta_0}(X_1^n)}{\int \pi(\theta) p_\theta(X_1^n) d\theta} \right].$$

Define the *Fisher score* (sum of gradients of log-likelihoods) $S_n := \sum_{i=1}^n \nabla \log p_{\theta_0}(X_i)$. For θ near θ_0 , for $X_i \stackrel{\text{iid}}{\sim} P_{\theta_0}$ we thus have

$$\begin{aligned} & \log p_\theta(X_1^n) \\ &= \sum_{i=1}^n \log p_{\theta_0}(X_i) + S_n^\top (\theta - \theta_0) + \frac{1}{2} \sum_{i=1}^n (\theta - \theta_0)^\top \nabla^2 \log p_{\theta_0}(X_i) (\theta - \theta_0) + o(\|\theta - \theta_0\|^2) \\ &\approx \log p_{\theta_0}(X_1^n) + S_n^\top (\theta - \theta_0) - \frac{n}{2} (\theta - \theta_0)^\top J_{\theta_0} (\theta - \theta_0), \end{aligned}$$

where we used the law of large numbers and the definition $J_\theta = -\mathbb{E}_{P_\theta}[\nabla^2 \log p_\theta(X)]$ of the Fisher information matrix, so that $\frac{1}{n} \sum_{i=1}^n \nabla^2 \log p_{\theta_0}(X_i) \rightarrow -J_{\theta_0}$ with probability 1. With this approximation, we complete the square so for $\bar{S}_n = \frac{1}{n} S_n$ we can write

$$\log p_\theta(X_1^n) \approx \log p_{\theta_0}(X_1^n) - \frac{n}{2} (\theta - \theta_0 - J_{\theta_0}^{-1} \bar{S}_n)^\top J_{\theta_0} (\theta - \theta_0 - J_{\theta_0}^{-1} \bar{S}_n) + \frac{n}{2} \bar{S}_n^\top J_{\theta_0}^{-1} \bar{S}_n \quad (19.5.6)$$

for θ near θ_0 .

We use the approximation (19.5.6) to compute expectations. First, observe that $\mathbb{E}_{\theta_0}[\bar{S}_n^\top J_{\theta_0}^{-1} \bar{S}_n] = \frac{d}{n}$, and so we typically have $J_{\theta_0}^{-1} \bar{S}_n = O(1/\sqrt{n})$. Thus, when θ is farther than distance order $1/\sqrt{n}$ from θ_0 , we have

$$\frac{n}{2} (\theta - \theta_0 - J_{\theta_0}^{-1} \bar{S}_n)^\top J_{\theta_0} (\theta - \theta_0 - J_{\theta_0}^{-1} \bar{S}_n) \rightarrow \infty.$$

Thus, assuming the approximation (19.5.6) is accurate enough, if we let Θ_0 be those points θ near θ_0 , then for “most” samples X_1^n we have

$$\begin{aligned} & \int \pi(\theta) p_\theta(X_1^n) d\theta \\ &\approx p_{\theta_0}(X_1^n) e^{\frac{n}{2} \bar{S}_n^\top J_{\theta_0}^{-1} \bar{S}_n} \int_{\Theta_0} \pi(\theta) \exp \left(-\frac{n}{2} (\theta - \theta_0 - J_{\theta_0}^{-1} \bar{S}_n)^\top J_{\theta_0} (\theta - \theta_0 - J_{\theta_0}^{-1} \bar{S}_n) \right) d\theta \\ &\approx p_{\theta_0}(X_1^n) \pi(\theta_0) e^{\frac{n}{2} \bar{S}_n^\top J_{\theta_0}^{-1} \bar{S}_n} \int_{\Theta_0} \exp \left(-\frac{n}{2} (\theta - \theta_0 - J_{\theta_0}^{-1} \bar{S}_n)^\top J_{\theta_0} (\theta - \theta_0 - J_{\theta_0}^{-1} \bar{S}_n) \right) d\theta \end{aligned}$$

because point far from θ_0 have essentially 0 mass in the integral and π is continuous. Treating the last integral as that of a Gaussian density (this is known as a *Laplace approximation*), where we recall that for $J \succ 0$ we have $\int_{\mathbb{R}^d} \exp(-x^\top J x) dx = (2\pi)^{d/2} \det(J)^{-1/2}$, we obtain

$$\int \pi(\theta) p_\theta(X_1^n) d\theta \approx p_{\theta_0}(X_1^n) \pi(\theta_0) \exp\left(\frac{n}{2} \bar{S}_n^\top J_{\theta_0}^{-1} \bar{S}_n\right) (2\pi)^{d/2} \det(nJ_{\theta_0})^{-1/2}. \quad (19.5.7)$$

Once we substitute the heuristic approximation (19.5.7) into the KL-divergence, we obtain

$$\begin{aligned} D_{\text{kl}}(P_{\theta_0}^n \| Q_\pi^n) &\approx \mathbb{E}_{P_{\theta_0}} \left[\log \frac{1}{\pi(\theta_0)} - \frac{n}{2} \bar{S}_n^\top J_{\theta_0}^{-1} \bar{S}_n - \frac{d}{2} \log(2\pi) + \frac{d}{2} \log n + \frac{1}{2} \log \det(J_{\theta_0}) \right] \\ &= \log \frac{1}{\pi(\theta_0)} + \frac{d}{2} \log \frac{1}{2\pi e} + \frac{d}{2} \log n + \frac{1}{2} \log \det(J_{\theta_0}), \end{aligned}$$

because $\mathbb{E}_{P_{\theta_0}}[\bar{S}_n \bar{S}_n^\top] = n^{-1} J_{\theta_0}$. Assuming each approximation term \approx adds an at most $o(1)$ error to the resulting formula, this is the result (19.5.4a). To make these statements rigorous requires arguments that control the error terms in the various integrals and expansions, typically via variants of Lebesgue's dominated convergence theorem; we view the heuristic (19.5.7) as satisfying.

19.6 Regret and capacity dualities

To develop rigorous versions of Theorem 19.3.2 requires some additional mathematical care. We will first prove the theorem in the special case that \mathcal{X} is finite, which corresponds to Corollary 19.3.4, because this contains the main ideas. The extension to arbitrary spaces \mathcal{X} essentially just requires more careful treatment of probabilistic convergence.

19.6.1 Duality when the domain is finite

To show that a regret/capacity duality obtains when \mathcal{X} is finite (and also when \mathcal{X} is infinite in the sections to come), we will leverage the Bregman divergence representation of expected regret (19.3.1). Let us recapitulate to keep the presentation self-contained. By Theorem 14.2.1, we have the Savage representation (14.2.1) associated to any proper loss,

$$\ell(Q, x) = -\Omega(Q) - \langle \nabla \Omega(Q), e_x - Q \rangle \quad (19.6.1)$$

for some convex Ω . In equation (19.6.1), we abuse notation slightly, in that we are (without comment) viewing distributions Q as elements of $\Delta_{\mathcal{X}} = \{q \in \mathbb{R}_+^{\mathcal{X}} \mid \langle \mathbf{1}, q \rangle = 1\} \subset \mathbb{R}^{\mathcal{X}}$, and e_x is the “ x th” standard basis vector. Then $\nabla \Omega(Q)$ is a (particular, fixed) element of the subdifferential $\partial \Omega(Q)$, which necessarily exists because Ω is proper (Theorem 14.2.1). As $\Omega(P) = \sup_Q -\mathbb{E}_P[\ell(Q, X)]$ for proper losses, we then have the equality (19.3.1), that is,

$$D_\Omega(P, Q) := \mathbb{E}_P[\ell(Q, X)] - \mathbb{E}_P[\ell(P, X)] = \Omega(P) - \Omega(Q) - \langle \nabla \Omega(Q), P - Q \rangle,$$

the Bregman divergence between P and Q . To avoid issues of non-determinacy in the selection of $\nabla \Omega(Q)$, we shall always take the gap in the expected losses as the definition of the divergence.

When the space \mathcal{X} is finite and the loss ℓ is strictly proper, the next lemma shows that convergence of distributions in D_Ω implies convergence in any norm. To make clear we work with vectors $p \in \mathbb{R}^{\mathcal{X}}$, we use lower case letters. We require one additional consideration: that

$$D_\Omega(p, q) := \mathbb{E}_p[\ell(q, X)] - \mathbb{E}_p[\ell(p, X)] = \mathbb{E}_p[\ell(q, X)] - H_\ell(p) = \Omega(p) + \mathbb{E}_p[\ell(q, X)]$$

is lower semicontinuous in p, q , meaning that if $p_n \rightarrow p$ and $q_n \rightarrow q$, then $\liminf_n D_\Omega(p_n, q_n) \geq D_\Omega(p, q)$. Because $\Omega(p)$ is the supremum of linear functions of p , it is lower semicontinuous, and so because \mathcal{X} is finite, it is enough that $\ell(q, x)$ be lower semicontinuous in q for each x .

Lemma 19.6.1. *Let \mathcal{X} be finite, and assume the loss ℓ is strictly proper and that D_Ω is lower semicontinuous. Then there exists a function ω with $\omega(\epsilon) > 0$ for all $\epsilon > 0$ such that*

$$D_\Omega(p, q) \geq \omega(\|p - q\|_1).$$

Proof We identify $\mathcal{P} = \Delta_{\mathcal{X}} = \{p \in \mathbb{R}_+^{\mathcal{X}} \mid \langle \mathbf{1}, p \rangle = 1\}$, a compact set. Assume for the sake of contradiction that

$$\inf_{p, q \in \Delta_{\mathcal{X}}} \{D_\Omega(p, q) \mid \|p - q\|_1 \geq \epsilon\} = 0$$

for some $\epsilon > 0$. Let $p_n, q_n \in \Delta_{\mathcal{X}}$ tend to this infimum, that is, satisfy $\|q_n - p_n\|_1 \geq \epsilon$ and $D_\Omega(p_n, q_n) \rightarrow 0$. Taking subsequences as necessary, we may assume that $q_n \rightarrow q$ and $p_n \rightarrow p$. By the lower semicontinuity assumption, for any $\delta > 0$ there exists N such that $n \geq N$ implies

$$\mathbb{E}_{p_n}[\ell(q_n, X)] - \mathbb{E}_{p_n}[\ell(p_n, X)] \geq \mathbb{E}_p[\ell(q, X)] + \Omega(p) - \delta = D_\Omega(p, q) - \delta.$$

Because $\delta > 0$ was arbitrary, we in fact have

$$\liminf_n \mathbb{E}_{p_n}[\ell(q_n, X)] - \mathbb{E}_{p_n}[\ell(p_n, X)] \geq D_\Omega(p, q) > 0$$

because ℓ is strictly proper, our desired contradiction. \square

As a consequence, if p_n is a Cauchy sequence for D_Ω , meaning that $\sup_{m \geq 0} D_\Omega(p_n, p_{n+m}) \rightarrow 0$ as $n \rightarrow \infty$, it is also a Cauchy sequence in $\mathbb{R}^{\mathcal{X}}$ and thus converges to some p .

19.6.2 Proof of Corollary 19.3.4

We turn now to the proof of Corollary 19.3.4. Recall the convex collection Π of prior distributions π on Θ . Let Q be any distribution on \mathcal{X} , (meaning we may identify it with a vector in $\mathbb{R}_+^{\mathcal{X}}$), and define the saddle objective

$$\begin{aligned} L(Q, \pi) &:= \int \mathbb{E}_{P_\theta}[\ell(q, X)] d\pi(\theta) - \int \mathbb{E}_{P_\theta}[\ell(P_\theta, X)] d\pi(\theta) \\ &= \int [\Omega(P_\theta) - \Omega(Q) - \langle \nabla \Omega(Q), P_\theta - Q \rangle] d\pi(\theta) = \int D_\Omega(P_\theta, Q) d\pi(\theta). \end{aligned}$$

We decompose the proof into several parts. We first recall from Lemma 19.3.1 that if $\hat{\pi}$ attains the supremum in the definition (19.3.4) of the information in the experiment $\{P_\theta\}$ for the loss ℓ , then $(P_{\hat{\pi}}, \hat{\pi})$ is a saddle point. Using this intermediate result, we use finite-dimensional approximations of the capacity $\sup_\pi I_\ell(\pi; X)$, which in turn give a sequence of (explicit) saddle points $Q_n = P_{\pi_n}$ minimizing the worst-case redundancy. Finally, with a bit of algebraic manipulation and the careful use of the relationship between Bregman divergences D_Ω and proper losses ℓ , we can use completeness to show that limits of the sequence Q_n necessarily exist, and they are unique.

Finite-dimensional approximation and a limiting Q With Lemma 19.3.1 in hand, we proceed via finite dimensional approximations. Let $\pi_n \in \Pi$ be a sequence of priors satisfying $I_\ell(\pi_n; X) \rightarrow \sup_{\pi \in \Pi} I_\ell(\pi; X) = C_\ell$. Then the set

$$\Pi_n := \text{Conv}\{\pi_1, \dots, \pi_n\} \quad (19.6.2)$$

is isomorphic to an at most $n - 1$ dimensional set, and so because the saddle objective $L(Q, \pi)$ is linear in π , $\inf_Q L(Q, \pi)$ is a concave function (and upper semicontinuous, as it is the infimum of linear functions). So

$$\pi_n^* = \operatorname{argmax}_{\pi \in \Pi_n} \left\{ \inf_Q L(Q, \pi) = I_\ell(\pi; X) \right\}$$

is attained. The sequence π_n^* of priors also satisfies $I_\ell(\pi_n^*; X) \rightarrow C_\ell$ because $C_\ell \geq I_\ell(\pi_n^*; X) \geq I_\ell(\pi_n; X) \rightarrow C_\ell$. Because the loss ℓ is proper, the mixture distribution $Q_n = P_{\pi_n^*}$ minimizes $L(Q, \pi_n^*)$, and Lemma 19.3.1 implies $\sup_{\pi \in \Pi_n} L(Q_n, \pi) \leq L(Q_n, \pi_n^*) \leq \inf_Q L(Q, \pi_n^*) \rightarrow C_\ell$.

Now for $k \in \mathbb{N}$ consider the pair $Q_n = P_{\pi_n^*}$ and $Q_{n+k} = P_{\pi_{n+k}^*}$. We have

$$L(Q_{n+k}, \pi_n^*) - L(Q_n, \pi_n^*) = -\Omega(Q_{n+k}) - \langle \nabla \Omega(Q_{n+k}), P_{\pi_n^*} - Q_{n+k} \rangle + \Omega(Q_n) = D_\Omega(Q_n, Q_{n+k}).$$

On the other hand, we also have

$$\begin{aligned} L(Q_{n+k}, \pi_n^*) - L(Q_n, \pi_n^*) &\leq L(Q_{n+k}, \pi_{n+k}^*) - L(Q_n, \pi_n^*) \\ &= I_\ell(\pi_{n+k}^*; X) - I_\ell(\pi_n^*; X) \rightarrow 0, \end{aligned}$$

the inequality following again from Lemma 19.3.1. As k is arbitrary, we have the limit (in n)

$$\sup_{k \geq 1} D_\Omega(Q_n, Q_{n+k}) \rightarrow 0. \quad (19.6.3)$$

Here we use that the space \mathcal{X} is finite; by Lemma 19.6.1, there is necessarily some $Q^* \in \Delta_{\mathcal{X}}$ such that $Q_n \rightarrow Q^*$. We show this Q^* is a saddle.

Optimality of the limiting Q Using Lemma 19.3.1 again, we have

$$L(Q_n, \pi_n^*) \leq I_\ell(\pi_n^*; X) \leq C_\ell.$$

For any $\pi \in \Pi_\infty := \cup_{n=1}^\infty \Pi_n(\Theta)$, we have $\pi \in \Pi_n(\Theta)$ for all large enough n , so that similarly

$$L(Q_n, \pi) \leq C_\ell$$

for all $\pi \in \Pi_\infty$. Now we may use the lower semi-continuity of the regret $D_\Omega(P_\theta, Q)$ in Q , which implies

$$\begin{aligned} \int (\mathbb{E}_{P_\theta}[\ell(Q^*, X)] - H_\ell(P_\theta)) d\pi(\theta) &\stackrel{(i)}{\leq} \int \liminf_n (\mathbb{E}_{P_\theta}[\ell(Q_n, X)] - H_\ell(P_\theta)) d\pi(\theta) \\ &\stackrel{(ii)}{\leq} \liminf_n \int (\mathbb{E}_{P_\theta}[\ell(Q_n, X)] - H_\ell(P_\theta)) d\pi(\theta) \leq C, \end{aligned}$$

where step (i) used lower semicontinuity and step (ii) used Fatou's lemma (Proposition A.3.1). In brief, we have demonstrated that

$$L(Q^*, \pi) \leq \sup_{\pi \in \Pi} I_\ell(\pi; X) = C_\ell \quad \text{for all } \pi \in \Pi_\infty = \bigcup_{n \geq 1} \Pi_n. \quad (19.6.4)$$

For the final step to showing that Q^* is optimal, we show that the inequality (19.6.4) holds for any $\pi \in \Pi$. To that end, let $\pi \in \Pi$ be any fixed prior, and then replace the set (19.6.2) with $\tilde{\Pi}_n = \text{Conv}\{\pi_1, \dots, \pi_n, \pi\} \subset \Pi$. This is isomorphic to an n -dimensional set, and we repeat precisely the same derivation *mutatis mutandis*, yielding a sequence of priors $\tilde{\pi}_n^*$ and distributions $\tilde{Q}_n = P_{\tilde{\pi}_n^*}$. Then

$$L(\tilde{Q}_n, \pi_n^*) - L(Q_n, \pi_n^*) = -\Omega(\tilde{Q}_n) - \langle \nabla \Omega(\tilde{Q}_n), P_{\pi_n^*} - \tilde{Q}_n \rangle + \Omega(Q_n) = D_\Omega(Q_n, \tilde{Q}_n).$$

Similarly, we have

$$\begin{aligned} L(\tilde{Q}_n, \pi_n^*) - L(Q_n, \pi_n^*) &\leq L(\tilde{Q}_n, \tilde{\pi}_n^*) - L(Q_n, \pi_n^*) \\ &= I_\ell(\tilde{\pi}_n^*; X) - I_\ell(\pi_n^*; X) \leq C_\ell - I_\ell(\pi_n^*; X) \rightarrow 0 \end{aligned}$$

by the assumption that $I_\ell(\pi_n; X) \rightarrow C_\ell = \sup_{\pi \in \Pi} I_\ell(\pi; X)$. Thus $D_\Omega(Q_n, \tilde{Q}_n) \rightarrow 0$, and so again by Lemma 19.6.1, $\tilde{Q}_n \rightarrow Q^*$ as well. Thus we extend inequality (19.6.4) to

$$L(Q^*, \pi) \leq \sup_{\pi \in \Pi} I_\ell(\pi; X) = C_\ell \text{ for all } \pi \in \Pi. \quad (19.6.5)$$

Saddle point and its uniqueness As an immediate consequence of inequality (19.6.5), we see that Q^* is indeed a saddle point: we have

$$\sup_{\pi \in \Pi} L(Q^*, \pi) = \sup_{\pi \in \Pi} I_\ell(\pi; X) = \sup_{\pi \in \Pi} \inf_Q L(Q, \pi) \leq \inf_Q \sup_{\pi \in \Pi} L(Q, \pi)$$

by the weak min-max inequality, so equality holds in each step. Finally, we show that Q^* is unique. Suppose that \tilde{Q} is a different distribution satisfying $C_\ell = \sup_{\pi \in \Pi} L(\tilde{Q}, \pi)$. Then letting

$$\begin{aligned} C_\ell &\geq L(\tilde{Q}, \pi_n^*) = L(\tilde{Q}, \pi_n^*) - L(Q_n, \pi_n^*) + L(Q_n, \pi_n^*) \\ &= -\Omega(\tilde{Q}) - \langle \nabla \Omega(\tilde{Q}), P_{\pi_n^*} - \tilde{Q} \rangle + \Omega(Q_n) + L(Q_n, \pi_n^*) = L(Q_n, \pi_n^*) + D_\Omega(Q_n, \tilde{Q}). \end{aligned}$$

But we know that $L(Q_n, \pi_n^*) \rightarrow C_\ell$, and thus

$$D_\Omega(Q_n, \tilde{Q}) \rightarrow 0. \quad (19.6.6)$$

Because $Q_n \rightarrow Q^*$, it must be the case that $\tilde{Q} = Q^*$.

If π^* achieves the capacity, that is, maximizes $\inf_Q L(Q, \pi)$ over $\pi \in \Pi$, then $\inf_Q L(Q, \pi^*) = C_\ell(\Pi)$, and by the inequality (19.6.5),

$$\sup_{\pi \in \Pi} L(Q^*, \pi) \leq \inf_Q L(Q, \pi^*).$$

Then $L(Q^*, \pi^*) \leq \inf_Q L(Q, \pi^*)$, and because ℓ is strictly proper we necessarily have $Q^* = P_{\pi^*}$.

19.6.3 Regret/capacity duality for arbitrary domains

Our proof of Corollary 19.3.4 in Section 19.6.2 made little use of the fact that \mathcal{X} was finite. An inspection of the proof shows that we used the finiteness of \mathcal{X} in only a few places: the first is in equation (19.6.1) representing the expected regret via proper losses and the Savage

representation (14.2.1) from Theorem 14.2.1. In this case, we can rewrite the loss using the general representation in Theorem 14.2.8,

$$\ell(Q, x) = -\Omega(Q) - \langle \nabla \Omega(Q), \mathbf{1}_x - Q \rangle \quad (19.6.7)$$

for some convex Ω . Here, we have abused notation to mimic that in (19.6.1): we take $\mathbf{1}_x$ to be the point mass at x and inner products to mean expectations, i.e.,

$$\langle \nabla \Omega(Q), P \rangle = \mathbb{E}_P[\Omega'(Q, X)]$$

for an appropriate subgradient $\Omega'(Q, x)$, as in expression (14.2.6). we then have the equality

$$D_\Omega(P, Q) := \mathbb{E}_P[\ell(Q, X)] - \mathbb{E}_P[\ell(P, X)] = \Omega(P) - \Omega(Q) - \langle \Omega(Q), P - Q \rangle,$$

exactly as in the finite-dimensional case. Therefore, as in the proof of Corollary 19.3.4, we still have the saddle point representation

$$L(Q, \pi) = \int D_\Omega(P_\theta, Q) d\pi(\theta)$$

and the infimal representation (19.3.6) that $\inf_Q L(Q, \pi) = L(P_\pi, \pi) = I_\ell(\pi; X)$, and so at least insofar as setting up the problem, finiteness of \mathcal{X} is immaterial because of the general proper loss representation in Theorem 14.2.8.

The other places in which we use finite dimensionality all revolve around continuity of divergence measures and in completeness of collections of distributions. The second place in which we use (something like) finite dimensionality is in the statement of Lemma 19.3.1, which assumes that the expected regret (Bregman divergence) $D_\Omega(P_\theta, Q)$ is lower semicontinuous along one-dimensional directions in Q . The remaining three uses of finite dimensionality all regard completeness: equation (19.6.3), uses that if Q_n is Cauchy for the Bregman divergence, i.e., $\sup_{k \geq 1} D_\Omega(Q_n, Q_{n+k}) \rightarrow 0$, then $Q_n \rightarrow Q^*$ for some distribution Q^* . We use this same completeness to argue that if $Q_n \rightarrow Q^*$ and $D_\Omega(Q_n, \tilde{Q}_n) \rightarrow 0$, then $\tilde{Q}_n \rightarrow Q^*$ as well in the derivation of inequality (19.6.5). And finally, the convergence (19.6.6) uses that $D_\Omega(Q_n, \tilde{Q}) \rightarrow 0$ and $D_\Omega(Q_n, Q) \rightarrow 0$ implies $Q = \tilde{Q}$, which proves that Q^* is unique.

Thus, while we will require somewhat careful topological and metric space considerations, to prove a rigorous version of Theorem 19.3.2 for arbitrary domains requires demonstrating only two conditions relating to modes of convergence of probability distributions:

- (A) The divergence $D_\Omega(P_\theta, Q)$ is directionally lower semicontinuous (19.3.8) in Q .
- (B) D_Ω is complete for the desired mode of convergence: if P_n is Cauchy for the divergence D_Ω , meaning that $\sup_{m \geq 1} D_\Omega(P_n, P_{n+m}) \rightarrow 0$, then there exists P such that $P_n \rightarrow P$, and if $D_\Omega(P_n, P) \rightarrow 0$ and $D_\Omega(P_n, Q) \rightarrow 0$, then $P = Q$.

If we can show both desiderata (A) and (B), then evidently Theorem 19.3.2 will be airtight. This requires a small detour into the convergence of distributions in different topologies, all of which coincide when the space \mathcal{X} is finite. (See Appendix A.4 for a discussion of different modes of convergence of probability distributions on metric spaces.) We shall focus on two modes of convergence: in total variation distance and convergence in distribution.

Convergence in total variation Convergence in total variation provides a more straightforward way to verify the above desiderata, thus extending Corollary 19.3.4 to arbitrary spaces. Let \mathcal{X} be an arbitrary space and P_n and P be probability measures on \mathcal{X} . Then we say that

$$P_n \rightarrow P \text{ in total variation if } \|P_n - P\|_{\text{TV}} \rightarrow 0.$$

Total variation is particularly convenient for convergence of distributions, as it is complete: if P_n is a Cauchy sequence for total variation, then it necessarily has a limit distribution P for which $\|P_n - P\|_{\text{TV}} \rightarrow 0$. (See Lemma A.4.4 in Appendix A.4.) Because we frequently work with divergences other than the variation distance, it is useful to allow other completeness notions.

Definition 19.1. Let \mathcal{P} be a collection of probability distributions on \mathcal{X} and $D : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_+$. Then D is complete for \mathcal{P} if

- i. Whenever $P_n \in \mathcal{P}$ is Cauchy for D , meaning $\sup_{k \in \mathbb{N}} D(P_n, P_{n+k}) \rightarrow 0$ as $n \rightarrow \infty$, then there exists $P \in \mathcal{P}$ for which $\|P_n - P\|_{\text{TV}} \rightarrow 0$
- ii. If $D(P_n, Q_n) \rightarrow 0$, then $\|P_n - Q_n\|_{\text{TV}} \rightarrow 0$.

For the KL-divergence, Pinsker's inequality essentially immediately provides completeness:

Example 19.6.2: Let $D(P, Q) = D_{\text{kl}}(P\|Q)$. Then Pinsker's inequality (Proposition 2.2.8) implies that $D(P, Q) \geq 2\|P - Q\|_{\text{TV}}^2$. Many choices for the set \mathcal{P} are possible here so long as \mathcal{P} is closed for the variation distance. As particular examples, if \mathcal{P} consists of the collection of all distributions on \mathcal{X} , then the KL-divergence is complete. Similarly, if \mathcal{P} consists of all distributions on \mathcal{X} absolutely continuous respect to some base measure ν , Definition 19.1 is satisfied as well (again, see Lemma A.4.4 in Appendix A.4).

The KL-divergence corresponds to the logarithmic loss $\ell(Q, x) = -\log q(x)$, as we have seen several times. Example 14.2.11 treats the case of general, potentially non-discrete, spaces \mathcal{X} . \diamond

Note that the divergence $D(P, Q) = D_{\text{kl}}(P\|Q)$ and saddle objective $L(Q, \pi) = \int D_{\text{kl}}(P_\theta\|Q) d\pi(\theta)$ are well-defined regardless of the domain \mathcal{X} , so that the redundancy/capacity theorem holds so long as we can verify the lower semicontinuity (A). For this, we have the following observation, which is a corollary of the Donsker-Varadhan representation of the KL-divergence (Theorem 6.1.1):

Corollary 19.6.3 (Semicontinuity of KL-divergence). Let $\|P_n - P\|_{\text{TV}}$ and $\|Q_n - Q\|_{\text{TV}} \rightarrow 0$. Then $\liminf_n D_{\text{kl}}(P_n\|Q_n) \geq D_{\text{kl}}(P\|Q)$.

Proof Let g be a simple function, that is, satisfying $g(x) = \sum_{i=1}^k \alpha_i \mathbf{1}\{x \in A_i\}$ for scalars α_i and measurable sets A_i . Then clearly $\mathbb{E}_{P_n}[g] \rightarrow \mathbb{E}_P[g]$ and $\mathbb{E}_{Q_n}[e^g] \rightarrow \mathbb{E}_Q[e^g]$ by the convergence in total variation. So

$$D_{\text{kl}}(P_n\|Q_n) \geq \mathbb{E}_{P_n}[g] - \log \mathbb{E}_{Q_n}[e^g] \rightarrow \mathbb{E}_P[g] - \log \mathbb{E}_Q[e^g].$$

As g was arbitrary, we may take a supremum over such simple functions to achieve the KL-divergence on the right. \square

Combining this corollary with the discussion in Example 19.6.2, we see we have satisfied desiderata (A) and (B), and have thus proved the redundancy/capacity duality in Corollary 19.3.3.

Convergence in distribution Let \mathcal{X} be a metric space (we leave the metric implicit). Recall that a sequence of distributions¹ P_n on \mathcal{X} *converge in distribution* to a distribution P , which we denote

$$P_n \xrightarrow{d} P,$$

if for every bounded continuous function $f : \mathcal{X} \rightarrow \mathbb{R}$,

$$\mathbb{E}_{P_n}[f(X)] \rightarrow \mathbb{E}_P[f(X)].$$

When $\mathcal{X} \subset \mathbb{R}^d$, this is equivalent to convergence of cumulative distribution functions at points of continuity: letting $F_n(t) = P_n(X \preceq t)$ and $F(t) = P(X \preceq t)$ be the CDFs of P_n and P , we have $P_n \xrightarrow{d} P$ if and only if

$$F_n(t) \rightarrow F(t)$$

at all continuity points of F . See Appendix A.4 for discussion around these modes of convergence.

Especially in infinite-dimensional problems, it can be easier to obtain completeness for convergence in distribution rather than in total variation, as we describe in Definition 19.1.

Definition 19.2. The divergence measure $D : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_+$ is complete for \mathcal{P} in distribution if

i. Whenever $P_n \in \mathcal{P}$ is Cauchy for D , there exists $P \in \mathcal{P}$ for which $P_n \xrightarrow{d} P$

ii. If $P_n \xrightarrow{d} P$ and $D(P_n, Q_n) \rightarrow 0$, then $Q_n \xrightarrow{d} P$.

To distinguish this from definition 19.1 and highlight the necessity of considering convergence in distribution, consider predicting cumulative distributions. As in forecasting problems, natural losses include the cumulative ranked probability score (CRPS) that we discuss in Examples 14.2.6 and 14.2.10.

Example 19.6.4: For distributions P and Q on \mathbb{R}^d , let

$$D(P, Q) = \int (P(X \preceq t) - Q(X \preceq t))^2 dt$$

be the average squared distance between their cumulative distribution functions (CDFs). Let \mathcal{P} be the collection of all distributions on a compact set $T \subset \mathbb{R}^d$. Then Proposition A.4.7 in Appendix A.4 shows the (more or less) standard result that this divergence is complete for convergence in distribution: we have $D(P_n, P) \rightarrow 0$ if and only if $P_n \xrightarrow{d} P$ and P_n is Cauchy for the Bregman divergence if and only if it has a limit $P_n \xrightarrow{d} P$. That is, it satisfies Definition 19.2.

In this case, the divergence corresponds to a (potentially) multivariate version of the CRPS loss $\ell_{\text{crps}}(Q, y) := \int (Q(Y \preceq t) - \mathbf{1}\{y \preceq t\})^2 dt$, which a minor generalization of Examples 14.2.6 and 14.2.10 makes clear. \diamond

Inspecting Definition (19.2) and Example 14.2.10, we see that both of the desiderata (A) and (B) hold: the semicontinuity follows because for $Q_n \xrightarrow{d} Q$,

$$\left| D(P, Q_n)^{1/2} - D(P, Q)^{1/2} \right| \leq D(Q, Q_n)^{1/2} \rightarrow 0$$

by the triangle inequality and the example. The completeness we have shown.

¹We shall assume all probabilities are Borel measures, meaning they are defined on the Borel sets.

19.6.4 A formal statement of regret/capacity duality

With the examples above in mind, we can provide a formal statement that makes Theorem 19.3.2 precise. All that we really need is an appropriate assumption on the losses ℓ to guarantee some type of completeness. Recall the representation (19.6.7) of a proper loss, so that for the negative generalized entropy Ω we have Bregman divergence

$$D_\Omega(P, Q) = \mathbb{E}_P[\ell(Q, X)] - \mathbb{E}_P[\ell(P, X)].$$

Definition 19.3. Let $\{P_\theta\}_{\theta \in \Theta}$ be a collection of distributions on \mathcal{X} indexed by $\theta \in \Theta$, let Π be a convex collection of distributions on Θ , and let $\mathcal{P} \supset \{P_\pi\}_{\pi \in \Pi}$. The loss ℓ is identifying for \mathcal{P} if

- i. it is strictly proper
- ii. for the negative generalized entropy Ω in its representation (19.6.7), either
 - a. the divergence D_Ω is complete for \mathcal{P} in total variation (Definition 19.1)
 - b. the divergence D_Ω is complete for \mathcal{P} in distribution (Definition 19.2)
- iii. for each $\theta \in \Theta$, $D_\Omega(P_\theta, Q)$ is directionally lower semicontinuous (19.3.8) in Q .

By our discussion so far, this is precisely what is needed to for a non-informal version of Theorem 19.3.2. Recalling the definition (19.3.5) of the saddle objective

$$L(Q, \pi) := \int (\mathbb{E}_{P_\theta}[\ell(Q, X)] - \mathbb{E}_{P_\theta}[\ell(P_\theta, X)]) d\pi(\theta),$$

we have the following theorem.

Theorem 19.6.5. Let $\{P_\theta\}_{\theta \in \Theta}$ be a collection of distributions on $X \in \mathcal{X}$, let Π be a convex collection of prior distributions on $\theta \in \Theta$. Let ℓ be identifying for $\mathcal{P} \supset \{P_\pi\}_{\pi \in \Pi}$ (Definition 19.3), and assume the capacity

$$C_\ell(\Pi) := \sup_{\pi \in \Pi} I_\ell(\pi; X) < \infty.$$

Then there exists a unique $Q^* \in \mathcal{P}$ such that

$$\sup_{\pi \in \Pi} L(Q^*, \pi) = \sup_{\pi \in \Pi} I_\ell(\pi; X),$$

If additionally $\Pi = \Pi(\Theta)$ consists of all distributions on Θ , then

$$\sup_{\theta \in \Theta} (\mathbb{E}_{P_\theta}[\ell(Q^*, X)] - \mathbb{E}_{P_\theta}[\ell(P_\theta, X)]) = C_\ell(\Pi(\Theta)) = \sup_{\pi \in \Pi} I_\ell(\pi; X).$$

Lastly, if π^* achieves the capacity $C_\ell(\Pi)$, then $Q^* = P_{\pi^*}$, and (P_{π^*}, π^*) is a saddle point:

$$\sup_{\pi \in \Pi} L(P_{\pi^*}, \pi) \leq L(P_{\pi^*}, \pi^*) \leq \inf_Q L(Q, \pi^*).$$

Clearly, Theorem 19.6.5 shows that we have the general regret/capacity duality (19.3.7).

19.7 Bibliographic details

See also the book of Grünwald [105] for more discussion of this and other issues. Gallager [97]. Discuss redundancy/capacity. My proof of Theorem 19.6.5 follows the same outline as the proof of Theorem 5.9 in Polyanskiy and Wu [155]. (See also Csiszár [60].) The generality of (arbitrary) proper losses requires a bit of care, but shows there is nothing particularly special that the log-loss buys.

19.8 Exercises

JCD Comment: A few exercise ideas: 1. General corollary for robust Bayes things, building off of Theorem 19.6.5. 2. Show how regret for log-loss means we get regret for other (bounded) losses by total variation. Originally had that in the book; can be just an exercise now.

JCD Comment: Do exercises with different entropies, e.g., $\Omega(p) = -\sum_x \sqrt{p(x)}$ or other power-type functions. These are Legendre.

Exercise 19.1 (Minimax redundancy and different loss functions): In this question, we consider expected losses under the Bernoulli distribution. Assume that $X_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$, meaning that $X_i = 1$ with probability p and $X_i = 0$ with probability $1 - p$. We consider four different loss functions, and their associated expected regret, for measuring the accuracy of our predictions of such X_i . For each of the four choices below, we prove expected regret bounds on

$$\text{Red}_n(\hat{\theta}, P, \ell) := \sum_{i=1}^n \mathbb{E}_P[\ell(\hat{\theta}(X_1^{i-1}), X_i)] - \inf_{\theta} \sum_{i=1}^n \mathbb{E}_P[\ell(\theta, X_i)], \quad (19.8.1)$$

where $\hat{\theta}$ is a predictor based on X_1, \dots, X_{i-1} at time i . Define $S_i = \sum_{j=1}^i X_j$ to be the partial sum up to time i . For each of parts (a)–(c), at time i use the predictor

$$\hat{\theta}_i = \hat{\theta}(X_1^{i-1}) = \frac{S_{i-1} + \frac{1}{2}}{i}.$$

- (a) Loss function: $\ell(\theta, x) = \frac{1}{2}(x - \theta)^2$. Show that $\text{Red}_n(\hat{\theta}, P, \ell) \leq C \cdot \log n$ where C is a constant.
- (b) Loss function: $\ell(\theta, x) = x \log \frac{1}{\theta} + (1 - x) \log \frac{1}{1-\theta}$, the usual log loss for predicting probabilities. Show that $\text{Red}_n(\hat{\theta}, P, \ell) \leq C \cdot \log n$ whenever the true probability $p \in (0, 1)$, where C is a constant. *Hint: Note that there exists a prior π for which $\hat{\theta}$ is a Bayes strategy. What is this prior?*
- (c) Loss function: $\ell(\theta, x) = |x - \theta|$. Show that $\text{Red}_n(\hat{\theta}, P, \ell) \geq c \cdot n$, where $c > 0$ is a constant, whenever the true probability $p \notin \{0, \frac{1}{2}, 1\}$.
- (d) **Extra credit:** Show that there is a numerical constant $c > 0$ such that for any procedure $\hat{\theta}$, the worst-case redundancy $\sup_{p \in [0, 1]} \text{Red}_n(\hat{\theta}, \text{Bernoulli}(p), \ell) \geq c\sqrt{n}$ for the absolute loss ℓ in part (c). Give a strategy attaining this redundancy.

Exercise 19.2: Fill in the details in the calculations for Example 19.5.3.

Exercise 19.3 (Strong versions of redundancy): Assume that for a given $\theta \in \Theta$ we draw $X_1^n \sim P_\theta$. We define the Bayes redundancy for a family of distributions $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$ as

$$C_n^\pi := \inf_Q \int D_{\text{kl}}(P_\theta \| Q) d\pi(\theta) = I_\pi(T; X_1^n),$$

where π is a probability measure on Θ , T is distributed according to π , and conditional on $T = \theta$, we draw $X_1^n \sim P_\theta$, and I_π denotes the mutual information when T is drawn according to π . Define the maximin redundancy $C_n^* := \sup_\pi C_n^\pi$ as the worst-case Bayes redundancy. We show that for “most” points θ under the prior π , if $\bar{Q} = \int P_\theta d\pi(\theta)$ is the mixture of all the P_θ under the prior π , then no distribution Q can have substantially better redundancy than \bar{Q} .

Consider any distribution Q on the set \mathcal{X} and let $\epsilon \in [0, 1]$, and define the set of points θ where Q is ϵ -better than the worst case redundancy as

$$B_\epsilon := \{\theta \in \Theta : D_{\text{kl}}(P_\theta \| Q) \leq (1 - \epsilon)C_n^*\}.$$

(a) Show that for any prior π , we have

$$\pi(B_\epsilon) \leq \frac{\log 2 + C_n^* - I_\pi(T; X_1^n)}{\epsilon C_n^*}.$$

As an aside, note this implies that if π_i is a sequence of priors tending to $\sup_\pi I_\pi(T; X_1^n)$ and the redundancy $C_n^* \rightarrow \infty$, then so long as $C_n^* - I_{\pi_i}(T; X_1^n) \ll \epsilon C_n^*$, we have $\pi_i(B_\epsilon) \approx 0$.

(b) Assume that π attains the supremum in the definition of C_n^* . Show that

$$\pi(B_\epsilon) \leq O(1) \cdot \exp(-\epsilon C_n^*).$$

Hint: Introduce the random variable Z to be 1 if the random variable $T \in B_\epsilon$ and 0 otherwise, then use that $Z \rightarrow T \rightarrow X_1^n$ forms a Markov chain, and expand the mutual information. For part (b), the inequality $\frac{1-x}{x} \log \frac{1}{1-x} \leq 1$ for all $x \in [0, 1]$ may be useful.

Exercise 19.4 (Mixtures are as good as point distributions): Let P be a Laplace(λ) distribution on \mathbb{R} , meaning that $X \sim P$ has density

$$p(x) = \frac{\lambda}{2} \exp(-\lambda|x|).$$

Assume that $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P$, and let P^n denote the n -fold product of P . In this problem, we compare the predictive performance of distributions from the normal location family $\mathcal{P} = \{N(\theta, \sigma^2) : \theta \in \mathbb{R}\}$ with the mixture distribution Q^π over \mathcal{P} defined by the normal prior distribution $N(\mu, \tau^2)$, that is, $\pi(\theta) = (2\pi\tau^2)^{-1/2} \exp(-(\theta - \mu)^2/2\tau^2)$.

(a) Let $P_{\theta, \Sigma}$ be the multivariate normal distribution with mean $\theta \in \mathbb{R}^n$ and covariance $\Sigma \in \mathbb{R}^{n \times n}$. What is $D_{\text{kl}}(P^n \| P_{\theta, \Sigma})$?

(b) Show that $\inf_{\theta \in \mathbb{R}^n} D_{\text{kl}}(P^n \| P_{\theta, \Sigma}) = D_{\text{kl}}(P^n \| P_{0, \Sigma})$, that is, the mean-zero normal distribution has the smallest KL-divergence from the Laplace distribution.

(c) Let Q_n^π be the mixture of the n -fold products in \mathcal{P} , that is, Q_n^π has density

$$q_n^\pi(x_1^n) = \int_{-\infty}^{\infty} \pi(\theta) p_\theta(x_1) \cdots p_\theta(x_n) d\theta,$$

where π is $N(0, \tau^2)$. What is $D_{\text{kl}}(P^n \| Q_n^\pi)$?

- (d) Show that the redundancy of Q_n^π under the distribution P is asymptotically nearly as good as the redundancy of any $P_\theta \in \mathcal{P}$, the normal location family (so P_θ has density $p_\theta(x) = (2\pi\sigma^2)^{-1/2} \exp(-(x - \theta)^2/2\sigma^2)$). That is, show that

$$\sup_{\theta \in \mathbb{R}} \mathbb{E}_P \left[\log \frac{1}{q_n^\pi(X_1^n)} - \log \frac{1}{p_\theta(X_1^n)} \right] = \mathcal{O}(\log n)$$

for any prior variance $\tau^2 > 0$ and any prior mean $\mu \in \mathbb{R}$, where the big-Oh hides terms dependent on τ^2, σ^2, μ^2 .

- (e) **Extra credit:** Can you give an interesting condition under which such redundancy guarantees hold more generally? That is, using Proposition 19.5.4 in the notes, give a general condition under which

$$\mathbb{E}_P \left[\log \frac{1}{q^\pi(X_1^n)} - \log \frac{1}{p_\theta(X_1^n)} \right] = o(n)$$

as $n \rightarrow \infty$, for all $\theta \in \Theta$.

Exercise 19.5: Prove Proposition 19.1.3.

Answer to 19.5: We have the equalizer condition that for $P \in \mathcal{P}_\alpha^{\text{lin}}$,

$$\mathbb{E}_P[-\log p_\theta(X)] = -\langle \alpha, \theta \rangle + A(\theta) = \mathbb{E}_{P_\theta}[-\log p_\theta(X)] = H_\nu(P_\theta),$$

where H is the Shannon entropy w.r.t. the measure ν . For $Q \ll \nu$, we know that Q has density q w.r.t. ν , so $\sup_{P \in \mathcal{P}_\alpha^{\text{lin}}} \mathbb{E}_P[-\log q(X)] \geq \mathbb{E}_{P_\theta}[\log \frac{1}{q(X)}] = D_{\text{kl}}(P_\theta \| Q) + H_\nu(P_\theta)$, which gives a gap unless $D_{\text{kl}}(P_\theta \| Q) = 0$.

Apply Theorem 14.4.7 to obtain that

$$\inf_{Q \ll \nu} \mathbb{E}_P[-\log q(X)] = \inf_{Q \ll \nu} D_{\text{kl}}(P \| Q) + \mathbb{E}_P[-\log p(X)] = H(P) \leq H(P_\theta)$$

for any $P \in \mathcal{P}_\alpha^{\text{lin}}$. So P_θ is indeed a saddle point as desired. \square

Part V

Appendices

Appendix A

Miscellaneous mathematical results

This appendix collects several mathematical results and some of the more advanced mathematical treatment required for full proofs of the results in the book. It is not a core part of the book, but it does provide readers who wish to see the measure-theoretic rigor necessary for some of our results, or otherwise, to dot the appropriate I's and cross the appropriate T's.

A.1 The roots of a polynomial

A.2 Measure-theoretic development of divergence measures

A.3 Integral convergence and completeness of probability spaces

In this section, we record several results on the convergence of integrals as well as divergence measures. These prove useful for, among other things, exchanging integration and differentiation or discussing completeness of probability spaces. We prove only those for which we do not know standard references.

The first two results are standard facts about integration and appear, for example, in Royden [163].

Proposition A.3.1 (Fatou's lemma). *Let f_n be a sequence of nonnegative measurable functions and let $f(x) = \liminf_n f_n(x)$, and let μ be a measure. Then f is measurable, and*

$$\int f d\mu \leq \liminf \int f_n d\mu.$$

Theorem A.3.2 (Dominated convergence). *Let $f_n, f : X \rightarrow Y$ where (X, μ) is a measure space and Y is a Banach space, where f_n are measurable. Assume that $f_n \rightarrow f$ either in μ -measure or for μ -almost every x , and that there exists an integrable dominating function g such that $g(x) \geq \|f_n(x)\|$ for all x and n , where $\int g d\mu < \infty$. Then*

$$\int f_n(x) d\mu(x) \rightarrow \int f(x) d\mu(x) \quad \text{and} \quad \int \|f_n(x) - f(x)\| d\mu(x) \rightarrow 0.$$

Proposition A.3.3 (An extended Scheffé's lemma). *Let $f_n : X \rightarrow Y$ where (X, μ) is a measure space and Y is a Banach space. Assume $f_n \rightarrow f$ either in μ measure or for μ -almost every x , and that $\limsup_n \int \|f_n(x)\|^p d\mu(x) \leq \int \|f(x)\|^p d\mu(x) < \infty$. Then $\int \|f_n(x) - f(x)\|^p d\mu(x) \rightarrow 0$.*

Proof By convexity of the norm, $\|f_n(x) - f(x)\|^p \leq 2^p \|f_n(x)\|^p + 2^p \|f(x)\|^p$, so

$$0 \leq 2^p \|f_n(x)\|^p + 2^p \|f(x)\|^p - \|f_n(x) - f(x)\|^p \rightarrow 2^{p+1} \|f(x)\|^p$$

μ -almost everywhere (or in μ -measure). Then Fatou's lemma that $\int h d\mu \leq \liminf_n \int h_n d\mu$ for nonnegative h_n converging to h in μ -measure implies

$$\begin{aligned} 2^{p+1} \int \|f(x)\|^p d\mu(x) &\leq \liminf_n \int (2^p \|f_n(x)\|^p + 2^p \|f(x)\|^p - \|f_n(x) - f(x)\|^p) d\mu(x) \\ &\leq 2^{p+1} \int \|f(x)\|^p d\mu(x) - \limsup_n \int \|f_n(x) - f(x)\|^p d\mu(x). \end{aligned}$$

Rearrange to obtain the limit. □

Corollary A.3.4. *If \mathcal{X} is discrete and $q_n \rightarrow q$ pointwise, $\|Q_n - Q\|_{TV} \rightarrow 0$.*

Proof Observe that $\sum_x q_n(x) = 1$ and $\sum_x q(x) = 1$, so Proposition A.3.3. implies $\sum_x |q_n(x) - q(x)| \rightarrow 0$. Of course, $\|Q_n - Q\|_{TV} = \frac{1}{2} \sum_x |q_n(x) - q(x)|$. □

A.4 Probabilistic convergence

In this appendix, we collect some of the background results on convergence of random variables and distributions we use.

A.4.1 Classical results on convergence in distribution

In Chapter 19, we sometimes use results on probability distributions converging in different models, including in distribution. Here we recapitulate a few of the main results in that direction. Recall that for a metric space (\mathcal{X}, ρ) , a sequence of probability distributions P_n on \mathcal{X} *converges in distribution* to P , written $P_n \xrightarrow{d} P$, if

$$\mathbb{E}_{P_n}[f(X)] \rightarrow \mathbb{E}_P[f(X)]$$

for all bounded continuous $f : \mathcal{X} \rightarrow \mathbb{R}$. There are many equivalent versions of convergence in distribution, and the Portmanteau theorem provides several characterizations. (See, e.g., Billingsley [30, Chapter 1.2] or van der Vaart and Wellner [186, Chapter 1.3].) Because we assume P_n and P are probability distributions, which induce random elements X_n and X , we do not need to address the measurability questions of van der Vaart and Wellner [186, Chapter 1.3].

Theorem A.4.1 (Portmanteau). *Let (\mathcal{X}, ρ) be a metric space and P_n, P be probability distributions on \mathcal{X} . The following are all equivalent to the convergence in distribution $P_n \xrightarrow{d} P$.*

- (i) $\mathbb{E}_{P_n}[f(X)] \rightarrow \mathbb{E}_P[f(X)]$ for all bounded, 1-Lipschitz continuous $f : \mathcal{X} \rightarrow \mathbb{R}$.
- (ii) For all upper semicontinuous $f : \mathcal{X} \rightarrow \mathbb{R}$ bounded from above, $\limsup_n \mathbb{E}_{P_n}[f(X)] \leq \mathbb{E}_P[f(X)]$.
- (iii) For all lower semicontinuous $f : \mathcal{X} \rightarrow \mathbb{R}$ bounded from below, $\liminf_n \mathbb{E}_{P_n}[f(X)] \geq \mathbb{E}_P[f(X)]$.

- (iv) For all open sets $O \subset \mathcal{X}$, $\liminf_n P_n(O) \geq P(O)$.
- (v) For all closed sets $C \subset \mathcal{X}$, $\limsup_n P_n(C) \leq P(C)$.
- (vi) For all continuity sets A of P , meaning sets for which the boundary $\text{bd } A = \text{cl } A \setminus \text{int } A$ satisfies $P(\text{bd } A) = 0$, $\lim_n P_n(A) = P(A)$.

If additionally $\mathcal{X} \subset \mathbb{R}^d$, then letting $F_n(t) = P_n(X \preceq t)$ and $F(t) = P(X \preceq t)$ denote the CDFs of P_n and P , then $P_n \xrightarrow{d} P$ is equivalent to

- (vii) $F_n(t) \rightarrow F(t)$ for all continuity points t of F .

The topology induced by convergence in distribution possesses several elegant properties, including compactness when the underlying space is compact. Prokhorov's theorem (see Billingsley [30, Chapter 1.5] or van der Vaart and Wellner [186, Chapter 1.3]) makes this clear. For the theorem, we require a family of probability measures defined on the same σ -algebra \mathcal{A} , meaning that $P \in \mathcal{P}$ are all mappings $P : \mathcal{A} \rightarrow [0, 1]$. Then \mathcal{P} is *tight* if for all ϵ , there is a compact K such that $P(K) \geq 1 - \epsilon$ for all $P \in \mathcal{P}$. The class \mathcal{P} is *sequentially compact* if for any sequence $P_n \in \mathcal{P}$, there is a subsequence $P_{n(m)}$ and probability measure Q on $(\mathcal{X}, \mathcal{A})$ for which $P_{n(m)} \xrightarrow{d} Q$ (though we need not have $Q \in \mathcal{P}$).

Theorem A.4.2 (Prokhorov). *Let \mathcal{P} be a collection of measures on $(\mathcal{X}, \mathcal{A})$.*

- (i) *If \mathcal{P} is tight, then it is sequentially compact.*
- (ii) *Assume that (\mathcal{X}, ρ) is separable and complete. Then if \mathcal{P} is sequentially compact, it is tight.*

When \mathcal{X} is separable and complete, we can also metrize convergence in distribution via the Lévy-Prokhorov distance: for a set $A \subset \mathcal{X}$, define its ϵ -enlargement $A^\epsilon := \{x \in \mathcal{X} \mid \rho(x, A) < \epsilon\}$, where $\rho(x, A) = \inf_{y \in A} \rho(x, y)$. Then for probability distributions P, Q , define

$$d_{\text{prob}}(P, Q) := \inf \{ \epsilon \geq 0 \mid P(A) \leq Q(A^\epsilon) + \epsilon \text{ and } Q(A) \leq P(A^\epsilon) + \epsilon \}.$$

This metrizes convergence in distribution on separable complete metric spaces [30, Theorem 6.8].

Theorem A.4.3. *Let \mathcal{X} be separable and complete. Then*

$$d_{\text{prob}}(P_n, P) \rightarrow 0 \text{ if and only if } P_n \xrightarrow{d} P.$$

Additionally, let \mathcal{P} be the collection of all Borel probability measures on \mathcal{X} . Then \mathcal{P} is complete and separable in the d_{prob} metric.

A.4.2 Assorted convergence results for probability distributions

In this appendix, we collect a few results that follow from the basic definitions of convergence of probability measures. We include them for completeness, because though they are more or less standard, we do not know specific references.

First, collections of probability distributions on a set are essentially always complete for the total variation distance.

Lemma A.4.4. *Let P_n be probability measures on a measurable space $(\mathcal{X}, \mathcal{F})$, where \mathcal{F} is a σ -field of subsets of \mathcal{X} . Let P_n be Cauchy for the variation distance. Then*

i. There exists a probability distribution P on \mathcal{X} for which $\|P_n - P\|_{\text{TV}} \rightarrow 0$.

ii. If additionally P_n are absolutely continuous with respect to a measure ν , then P is as well.

Proof We show claim i first. For each $A \in \mathcal{F}$, we have $P_n(A) \in [0, 1]$, and $P_n(A)$ is Cauchy in \mathbb{R} and so has a limit we shall call $P(A)$ (though we have not demonstrated it is a probability distribution). It is immediate that this limit satisfies $P(\emptyset) = 0$ and $P(\mathcal{X}) = 1$, and that it is finitely additive. It remains to demonstrate that P is countably additive. For this, let A_1, A_2, \dots be a collection of disjoint measurable sets, and let $\epsilon > 0$ be arbitrary. For $m \in \mathbb{N}$, define the tail sets $B_m = \cup_{i>m} A_i$. For all $m \in \mathbb{N}$, finite additivity implies

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^m P(A_i) + P(B_m).$$

Now, there is some N such that $n \geq N$ implies

$$\sup_k \sup_m |P_n(B_m) - P_{n+k}(B_m)| \leq \sup_k \|P_n - P_{n+k}\|_{\text{TV}} \leq \epsilon$$

because P_n is Cauchy for the total variation. Taking $k \rightarrow \infty$, we obtain that if $n \geq N$,

$$|P_n(B_m) - P(B_m)| \leq \epsilon \text{ for all } m \in \mathbb{N}.$$

For any fixed n , the probability of the tail sets $P_n(B_m) \rightarrow 0$ as $m \rightarrow \infty$. Take m large enough that $P_n(B_m) \leq \epsilon$. Then $0 \leq P(B_m) \leq 2\epsilon$, so P satisfies

$$\sum_{i=1}^m P(A_i) \leq \sum_{i=1}^m P(A_i) + P(B_m) = P\left(\bigcup_{i=1}^{\infty} A_i\right) \leq \sum_{i=1}^m P(A_i) + 2\epsilon$$

for large m . As $\epsilon > 0$ was arbitrary, taking $m \rightarrow \infty$ gives that P is indeed countably additive.

It remains to show that $\|P_n - P\|_{\text{TV}} \rightarrow 0$. Fix $\epsilon > 0$ and let N be large enough that $\|P_n - P_{n+k}\|_{\text{TV}} \leq \epsilon$ for all $n \geq N$ and $k \geq 0$. For any set A , we have $P_{n+k}(A) \rightarrow P(A)$, so

$$|P_n(A) - P(A)| \leq \epsilon$$

for all $n \geq N$, which implies the result.

For the final claim of the lemma, take a set A such that $\nu(A) = 0$. Then $P_n(A) = 0$ for all n , and thus $P(A) = |P(A) - P_n(A)| \rightarrow 0$, i.e., $P(A) = 0$. \square

We now return to alternative metrics for convergence in distribution. For probability distributions P and Q on \mathbb{R}^d , we define the squared L^2 -distance of their cumulative distributions

$$D_{\text{cdf}}(P, Q) := \int (P(X \preceq t) - Q(X \preceq t))^2 dt.$$

Because the CDF $t \mapsto P(X \preceq t)$ is non-decreasing, we expect that the points of discontinuity should somehow be small, and thus we might hope that D_{cdf} metrizes convergence in distribution. Over compact sets, this holds. To demonstrate it fully rigorously, however, we present two results on the sets of discontinuities of functions.

Lemma A.4.5. Let $f : X \rightarrow Y$ be any function between metric spaces (X, ρ_X) and (Y, ρ_Y) , and let $\mathbb{B}_\delta(x) = \{z \in X \mid \rho_X(x, z) < \delta\}$ be the open δ -ball around x . Define the set

$$D_\epsilon(f) := \{x \in X \mid \text{for all } \delta > 0 \text{ there exist } y, z \in \mathbb{B}_\delta(x) \text{ s.t. } d_Y(f(y), f(z)) \geq \epsilon\}$$

of ϵ -discontinuities of f . Then $D_\epsilon(f)$ is closed, and the set

$$D(f) := \{x \in X \mid f \text{ is discontinuous at } x\}$$

is Borel-measurable.

Proof We show that $D_\epsilon(f)$ is closed. Let x_∞ be a limit point of $D_\epsilon(f)$. If x_∞ is isolated, meaning that there is some $\delta_0 > 0$ such that $(\mathbb{B}_{\delta_0}(x_\infty) \setminus \{x_\infty\}) \cap D_\epsilon(f) = \emptyset$, then clearly $x_\infty \in D_\epsilon(f)$. Otherwise, let $\delta > 0$ be arbitrary. Then there is some $x \in D_\epsilon(f)$ with $x \in \mathbb{B}_\delta(x_\infty) \setminus \{x_\infty\}$ by definition of the limit point, and because $\mathbb{B}_\delta(x_\infty) \setminus \{x_\infty\}$ is open, there exists $\delta_0 > 0$ such that $\mathbb{B}_{\delta_0}(x) \subset \mathbb{B}_\delta(x_\infty) \setminus \{x_\infty\}$. As $x \in D_\epsilon(f)$, there are certainly $y, z \in \mathbb{B}_{\delta_0}(x)$ such that $\rho_Y(f(y), f(z)) \geq \epsilon$. But then $y, z \in \mathbb{B}_\delta(x_\infty)$, and as δ was arbitrary we have $x_\infty \in D_\epsilon(f)$.

The Borel measurability of $D(f)$ follows because $D(f) = \bigcup_{n \geq 1} D_{1/n}(f)$ is the countable union of closed sets. \square

Lemma A.4.6. Let F be the cumulative distribution function of $X \in \mathbb{R}^d$. If $d = 1$, the set $D(F)$ of discontinuities of F is countable, and if $d > 1$, it is Borel measurable and has measure 0.

Proof The set of discontinuity points of a one-dimensional CDF F is necessarily countable: for $\epsilon > 0$, the set $D_\epsilon(F) := \{t \in \mathbb{R} \mid F(t) \geq \limsup_{\delta \downarrow 0} F(t - \delta) + \epsilon\}$ has cardinality at most $1/\epsilon < \infty$, so

$$D(F) := \left\{ t \in \mathbb{R} \mid \lim_{\delta \downarrow 0} F(t - \delta) < F(t) \right\} = \bigcup_{n \in \mathbb{N}} D_{1/n}(F)$$

is countable.

Let F be the cumulative distribution function of $X \in \mathbb{R}^d$ and F_1, \dots, F_d be the (marginal) cumulative distributions of X_1, \dots, X_d . Let $t, t' \in \mathbb{R}^d$, and define the elementwise maximum $t \vee t' = [\max\{t_j, t'_j\}]_{j=1}^d$ and similarly the elementwise minimum $t \wedge t'$. Then by monotonicity of the CDFs,

$$\begin{aligned} |F(t) - F(t')| &\leq F(t \vee t') - F(t \wedge t') = \mathbb{P}(t \wedge t' \prec X \preceq t \vee t') \\ &\leq \sum_{j=1}^d \mathbb{P}(\min\{t_j, t'_j\} < X_j \leq \max\{t_j, t'_j\}) = \sum_{j=1}^d (F_j(\max\{t_j, t'_j\}) - F_j(\min\{t_j, t'_j\})). \end{aligned}$$

so a point t is a discontinuity of F only if it is a discontinuity for at least one of the marginal CDFs F_j . That is, we have shown that the set $D(F)$ of discontinuities of F satisfies

$$D(F) \subset \bigcup_{j=1}^d \left\{ \mathbb{R}^{j-1} \times D(F_j) \times \mathbb{R}^{d-j} \right\}$$

$D(F)$ is Borel measurable by Lemma A.4.5. If λ denotes d -dimensional Lebesgue measure, monotonicity of λ then gives

$$\lambda(D(F)) \leq \sum_{j=1}^d \lambda\left(\mathbb{R}^{j-1} \times D(F_j) \times \mathbb{R}^{d-j}\right) \leq \sum_{j=1}^d \sum_{t \in D(F_j)} \lambda\left(\mathbb{R}^{j-1} \times \{t\} \times \mathbb{R}^{d-j}\right) = 0,$$

because the Lebesgue measure of any lower-dimensional subset of \mathbb{R}^d is 0. Lebesgue and Borel measures agree on the Borel sets. \square

Proposition A.4.7. *For measures P, Q on \mathbb{R}^d define*

$$D_{\text{cdf}}(P, Q) = \int (P(X \preceq t) - Q(X \preceq t))^2 dt.$$

Let T be a compact subset of \mathbb{R}^d . Then D_{cdf} is complete for convergence in distribution over T , and $P_n \xrightarrow{d} P$ if and only if $D_{\text{cdf}}(P_n, P) \rightarrow 0$.

Proof We first show that if $P_n \xrightarrow{d} P$, then $D_{\text{cdf}}(P_n, P) \rightarrow 0$. To see this, recall from Lemma A.4.6 that the points of discontinuity of any cumulative distribution function $t \mapsto P(X \preceq t)$ have measure 0 (and are also measurable). Thus, by definition of convergence in distribution, the (Lebesgue) measure of t for which $P_n(X \preceq t) \not\rightarrow P(X \preceq t)$ is 0. Now, we have a sequence of functions $F_n(t) := P_n(X \preceq t)$ and $F(t) := P(X \preceq t)$, where $F_n(t) \rightarrow F(t)$ for almost all t . Then because $(F_n(t) - F(t))^2 \leq \mathbf{1}\{t \in T\}$ and T is compact (so that $\int \mathbf{1}\{t \in T\} dt = \text{Vol}(T) < \infty$), Lebesgue's dominated convergence theorem implies $D_{\text{cdf}}(P_n, P) = \int (F_n(t) - F(t))^2 dt \rightarrow 0$.

Now we show that if $D_{\text{cdf}}(P_n, P) \rightarrow 0$, then $P_n \xrightarrow{d} P$. By Prokhorov's theorem (Thm. A.4.2), any subsequence $P_{n(m)}$ has a further subsequence that converges to some limit distribution; call this Q . Then by the triangle inequality,

$$D_{\text{cdf}}(P, Q)^{1/2} \leq D_{\text{cdf}}(P, P_{n(m)})^{1/2} + D_{\text{cdf}}(P_{n(m)}, Q)^{1/2} \rightarrow 0,$$

so $Q = P$. Using the standard topological result that if for a sequence P_n , every subsequence has a further subsequence converging to P , then P_n converges to P , we have the result.

Lastly, we show completeness. Let P_n be a Cauchy sequence for D_{cdf} . Then by Prokhorov's theorem, because T is compact, for every subsequence of P_n there is a further subsequence for which $P_{n(m)} \xrightarrow{d} Q$ for some probability distribution Q . Let Q_0, Q_1 be two subsequential limits; we show that $Q_0 = Q_1$. To see this, note that for each $n, m \in \mathbb{N}$,

$$D_{\text{cdf}}(Q_0, Q_1)^{1/2} \leq D_{\text{cdf}}(Q_0, P_n)^{1/2} + D_{\text{cdf}}(P_n, P_m)^{1/2} + D_{\text{cdf}}(P_m, Q_1)^{1/2}.$$

Now, choose n and m along subsequences that, respectively, satisfy $P_n \xrightarrow{d} Q_0$ and $P_m \xrightarrow{d} Q_1$. Then the first and last terms above converge to 0 by the first part of the proof, and by assumption that P_n is Cauchy, the middle term converges to 0. So $D_{\text{cdf}}(Q_0, Q_1) = 0$, and $Q_0 = Q_1$. \square

A.5 Stirling approximations and entropy

We develop several approximations to factorials and binomials here.

Lemma A.5.1 (Weak Stirling approximation). *Let $n \in \mathbb{N}$. Then*

$$\left(\frac{n}{e}\right)^n \leq n! \leq \min\{n, e^4/4\} \left(\frac{n}{e}\right)^n.$$

Proof For the upper bound,

$$\log n! = \sum_{i=2}^{n-1} \log i + \log n \leq \log n + \int_2^{n-1} \log x \, dx,$$

and as $\frac{\partial}{\partial x}(x \log x - x) = \log x$, we obtain

$$\begin{aligned} \log n! &\leq \log n + (x \log x - x)|_{x=2}^{n-1} = \log n + (n-1) \log(n-1) - (n-1) \log e - 2 \log 2 + 2 \\ &= \log n + (n-1) \log \frac{n-1}{e} + (3 - 2 \log 2). \end{aligned}$$

That is,

$$n! \leq n \left(\frac{n-1}{e} \right)^{n-1} \frac{e^3}{4} = n \left(\frac{n}{e} \right)^n \frac{e^4}{4n} \left(\frac{n-1}{n} \right)^{n-1} \leq \frac{e^4}{4} \left(\frac{n}{e} \right)^n.$$

We can explicitly check the inequality for $n \leq 13$.

For the lower bound, we observe that

$$\log n! = \sum_{i=1}^n \log i \geq \int_0^n \log x \, dx = (x \log x - x)|_{x=0}^n = n \log n - n \log e,$$

which is the desired result. □

Appendix B

Convex Analysis

In this appendix, we review several results in convex analysis that are useful for our purposes. We give only a cursory study here, identifying the basic results and those that will be of most use to us; the field of convex analysis as a whole is vast. The study of convex analysis and optimization has become very important practically in the last forty to fifty years for a few reasons, the most important of which is probably that convex optimization problems—those optimization problems in which the objective and constraints are convex—are tractable, while many others are not. We do not focus on optimization ideas here, however, building only some analytic tools that we will find useful. We borrow most of our results from Hiriart-Urruty and Lemaréchal [111], focusing mostly on the finite-dimensional case (though we present results that apply in infinite dimensional cases with proofs that extend straightforwardly, and we do not specify the domains of our functions unless necessary), as we require no results from infinite-dimensional analysis.

In addition, we abuse notation and assume that the range of any function is the *extended real line*, meaning that if $f : C \rightarrow \mathbb{R}$ we mean that $f(x) \in \mathbb{R} \cup \{-\infty, +\infty\}$, where $-\infty$ and $+\infty$ are infinite and satisfy $a + \infty = +\infty$ and $a - \infty = -\infty$ for any $a \in \mathbb{R}$. However, we assume throughout and without further mention that our functions are *proper*, meaning that $f(x) > -\infty$ for all x , as this allows us to avoid annoying pathologies.

B.1 Convex sets

We begin with the simplest and most important object in convex analysis, a convex set.

Definition B.1. A set C is convex if for all $\lambda \in [0, 1]$ and all $x, y \in C$, we have

$$\lambda x + (1 - \lambda)y \in C.$$

An important restriction of convex sets is to *closed* convex sets, those convex sets that are, well, closed.

JCD Comment: Picture

We now consider two operations that extend sets, convexifying them in nice ways.

Definition B.2. The affine hull of a set C is the smallest affine set containing C . That is,

$$\text{aff}(C) := \left\{ \sum_{i=1}^k \lambda_i x_i : k \in \mathbb{N}, x_i \in C, \lambda \in \mathbb{R}^k, \sum_{i=1}^k \lambda_i = 1 \right\}.$$

Associated with any set is also its convex hull:

Definition B.3. The convex hull of a set $C \subset \mathbb{R}^d$, denoted $\text{Conv}(C)$, is the intersection of all convex sets containing C .

JCD Comment: picture

An almost immediate associated result is that the convex hull of a set is equal to the set of all convex combinations of points in the set.

Proposition B.1.1. Let C be an arbitrary set. Then

$$\text{Conv}(C) = \left\{ \sum_{i=1}^k \lambda_i x_i : k \in \mathbb{N}, x_i \in C, \lambda \in \mathbb{R}_+^k, \sum_{i=1}^k \lambda_i = 1 \right\}.$$

Proof Call T the set on the right hand side of the equality in the proposition. Then $T \supset C$ is clear, as we may simply take $\lambda_1 = 1$ and vary $x \in C$. Moreover, the set $T \subset \text{Conv}(C)$, as any convex set containing C must contain all convex combinations of its elements; similarly, any convex set $S \supset C$ must have $S \supset T$.

Thus if we show that T is convex, then we are done. Take any two points $x, y \in T$. Then $x = \sum_{i=1}^k \alpha_i x_i$ and $y = \sum_{i=1}^l \beta_i y_i$ for $x_i, y_i \in C$. Fix $\lambda \in [0, 1]$. Then $(1 - \lambda)\beta_i \geq 0$ and $\lambda\alpha_i \geq 0$ for all i ,

$$\lambda \sum_{i=1}^k \alpha_i + (1 - \lambda) \sum_{i=1}^l \beta_i = \lambda + (1 - \lambda) = 1,$$

and $\lambda x + (1 - \lambda)y$ is a convex combination of the points x_i and y_i weighted by $\lambda\alpha_i$ and $(1 - \lambda)\beta_i$, respectively. So $\lambda x + (1 - \lambda)y \in T$ and T is convex. \square

We also give one more definition, which is useful for dealing with some pathological cases in convex analysis, as it allows us to assume many sets are full-dimensional.

Definition B.4. The relative interior of a set C is the interior of C relative to its affine hull, that is,

$$\text{relint}(C) := \{x \in C : B(x, \epsilon) \cap \text{aff}(C) \subset C \text{ for some } \epsilon > 0\},$$

where $B(x, \epsilon) := \{y : \|y - x\| < \epsilon\}$ denotes the open ball of radius ϵ centered at x .

An example may make Definition B.4 clearer.

Example B.1.2 (Relative interior of a disc): Consider the (convex) set

$$C = \left\{ x \in \mathbb{R}^d : x_1^2 + x_2^2 \leq 1, x_j = 0 \text{ for } j \in \{3, \dots, d\} \right\}.$$

The affine hull $\text{aff}(C) = \mathbb{R}^2 \times \{0\} = \{(x_1, x_2, 0, \dots, 0) : x_1, x_2 \in \mathbb{R}\}$ is simply the (x_1, x_2) -plane in \mathbb{R}^d , while the relative interior $\text{relint}(C) = \{x \in \mathbb{R}^d : x_1^2 + x_2^2 < 1\} \cap \text{aff}(C)$ is the “interior” of the 2-dimensional disc in \mathbb{R}^d . \diamond

In finite dimensions, we may actually restrict the definition of the convex hull of a set C to convex combinations of a bounded number (the dimension plus one) of the points in C , rather than arbitrary convex combinations as required by Proposition B.1.1. This result is known as *Carathéodory’s theorem*.

Theorem B.1.3. *Let $C \subset \mathbb{R}^d$. Then $x \in \text{Conv}(C)$ if and only if there exist points $x_1, \dots, x_{d+1} \in C$ and $\lambda \in \mathbb{R}_+^{d+1}$ with $\sum_{i=1}^{d+1} \lambda_i = 1$ such that*

$$x = \sum_{i=1}^{d+1} \lambda_i x_i.$$

Proof It is clear that if x can be represented as such a sum, then $x \in \text{Conv}(C)$. Conversely, Proposition B.1.1 implies that for any $x \in \text{Conv}(C)$ we have

$$x = \sum_{i=1}^k \lambda_i x_i, \quad \lambda_i \geq 0, \quad \sum_{i=1}^k \lambda_i = 1, \quad x_i \in C$$

for some λ_i, x_i . Assume that $k > d+1$ and $\lambda_i > 0$ for each i , as otherwise, there is nothing to prove. Then we know that the points $x_i - x_1$ are certainly linearly dependent (as there are $k-1 > d$ of them), and we can find (not identically zero) values $\alpha_2, \dots, \alpha_k$ such that $\sum_{i=2}^k \alpha_i (x_i - x_1) = 0$. Let $\alpha_1 = -\sum_{i=2}^k \alpha_i$ to obtain that we have both

$$\sum_{i=1}^k \alpha_i x_i = 0 \quad \text{and} \quad \sum_{i=1}^k \alpha_i = 0. \quad (\text{B.1.1})$$

Notably, the equalities (B.1.1) imply that at least one $\alpha_i > 0$, and if we define $\lambda^* = \min_{i: \alpha_i > 0} \frac{\lambda_i}{\alpha_i} > 0$, then setting $\lambda'_i = \lambda_i - \lambda^* \alpha_i$ we have

$$\lambda'_i \geq 0 \text{ for all } i, \quad \sum_{i=1}^k \lambda'_i = \sum_{i=1}^k \lambda_i - \lambda^* \sum_{i=1}^k \alpha_i = 1, \quad \text{and} \quad \sum_{i=1}^k \lambda'_i x_i = \sum_{i=1}^k \lambda_i x_i - \lambda^* \sum_{i=1}^k \alpha_i x_i = x.$$

But we know that at least one of the $\lambda'_i = 0$, so that we could write x as a convex combination of $k-1$ elements. Repeating this strategy until $k = d+1$ gives the theorem. \square

B.1.1 Operations preserving convexity

We now touch on a few simple results about operations that preserve convexity of convex sets. First, we make the following simple observation.

Observation B.1.4. *Let C be a convex set. Then $C = \text{Conv}(C)$.*

Observation B.1.4 is clear, as we have $C \subset \text{Conv}(C)$, while any other convex $S \supset C$ clearly satisfies $S \supset \text{Conv}(C)$. Secondly, we note that intersections preserve convexity.

Observation B.1.5. *Let $\{C_\alpha\}_{\alpha \in \mathcal{A}}$ be an arbitrary collection of convex sets. Then*

$$C = \bigcap_{\alpha \in \mathcal{A}} C_\alpha$$

is convex. Moreover, if C_α is closed for each α , then C is closed as well.

The convexity property follows because if $x_1 \in C$ and $x_2 \in C$, then clearly $x_1, x_2 \in C_\alpha$ for all $\alpha \in \mathcal{A}$, and moreover $\lambda x_1 + (1 - \lambda)x_2 \in C_\alpha$ for all α and any $\lambda \in [0, 1]$. The closure property is standard. In addition, we note that closing a convex set maintains convexity.

Observation B.1.6. *Let C be convex. Then $\text{cl}(C)$ is convex.*

To see this, we note that if $x, y \in \text{cl}(C)$ and $x_n \rightarrow x$ and $y_n \rightarrow y$ (where $x_n, y_n \in C$), then for any $\lambda \in [0, 1]$, we have $\lambda x_n + (1 - \lambda)y_n \in C$ and $\lambda x_n + (1 - \lambda)y_n \rightarrow \lambda x + (1 - \lambda)y$. Thus we have $\lambda x + (1 - \lambda)y \in \text{cl}(C)$ as desired.

Observation B.1.6 also implies the following result.

Observation B.1.7. *Let D be an arbitrary set. Then*

$$\bigcap \{C \text{ closed} : C \supset D, C \text{ is convex}\} = \text{cl Conv}(D).$$

Proof Let T denote the leftmost set. It is clear that $T \subset \text{cl Conv}(D)$ as $\text{cl Conv}(D)$ is a closed convex set (by Observation B.1.6) containing D . On the other hand, if $C \supset D$ is a closed convex set, then $C \supset \text{Conv}(D)$, while the closedness of C implies it also contains the closure of $\text{Conv}(D)$. Thus $T \supset \text{cl Conv}(D)$ as well. \square

JCD Comment: Picture

Observation B.1.8. *Let $C \subset \mathbb{R}^d$ be compact. Then $\text{Conv}(C)$ is compact.*

Proof Let $x^n \in \text{Conv}(C)$ converge to some x . By Theorem B.1.3, for each n we can find $x_1^n, \dots, x_{d+1}^n \in C$ and $\lambda^n \in \Delta_{d+1}$ such that $x^n = \sum_{j=1}^{d+1} \lambda_j^n x_j^n$. Taking subsequences as necessary, we can assume that $x_1^n \rightarrow x_1 \in C$, as C is compact. Then taking a further subsequence, we can as well assume $x_2^n \rightarrow x_2 \in C$, and so on, so that $x_j^n \rightarrow x_j$ for each $j = 1, \dots, d+1$. Then because Δ_{d+1} is compact, we can likewise take a (further) subsequence of λ^n so that $\lambda^n \rightarrow \lambda \in \Delta_{d+1}$. Evidently these limiting objects satisfy $x = \sum_{j=1}^{d+1} \lambda_j x_j$. \square

As our last consideration of operations that preserve convexity, we consider what is known as the perspective of a set. To define this set, we need to define the perspective function, which, given a point $(x, t) \in \mathbb{R}^d \times \mathbb{R}_{++}$ (here $\mathbb{R}_{++} = \{t : t > 0\}$ denotes strictly positive points), is defined as

$$\text{pers}(x, t) = \frac{x}{t}.$$

We have the following definition.

Definition B.5. *Let $C \subset \mathbb{R}^d \times \mathbb{R}_+$ be a set. The perspective transform of C , denoted by $\text{pers}(C)$, is*

$$\text{pers}(C) := \left\{ \frac{x}{t} : (x, t) \in C \text{ and } t > 0 \right\}.$$

This corresponds to taking all the points $z \in C$, normalizing them so their last coordinate is 1, and then removing the last coordinate. (For more on perspective functions, see Boyd and Vandenberghe [38, Chapter 2.3.3].)

It is interesting to note that the perspective of a convex set is convex. First, we note the following.

Lemma B.1.9. *Let $C \subset \mathbb{R}^{d+1}$ be a compact line segment, meaning that $C = \{\lambda x + (1 - \lambda)y : \lambda \in [0, 1]\}$, where $x_{d+1} > 0$ and $y_{d+1} > 0$. Then $\text{pers}(C) = \{\lambda \text{pers}(x) + (1 - \lambda) \text{pers}(y) : \lambda \in [0, 1]\}$.*

Proof Let $\lambda \in [0, 1]$. Then

$$\begin{aligned} \text{pers}(\lambda x + (1 - \lambda)y) &= \frac{\lambda x_{1:d} + (1 - \lambda)y_{1:d}}{\lambda x_{d+1} + (1 - \lambda)y_{d+1}} \\ &= \frac{\lambda x_{d+1}}{\lambda x_{d+1} + (1 - \lambda)y_{d+1}} \frac{x_{1:d}}{x_{d+1}} + \frac{(1 - \lambda)y_{d+1}}{\lambda x_{d+1} + (1 - \lambda)y_{d+1}} \frac{y_{1:d}}{y_{d+1}} \\ &= \theta \text{pers}(x) + (1 - \theta) \text{pers}(y), \end{aligned}$$

where $x_{1:d}$ and $y_{1:d}$ denote the vectors of the first d components of x and y , respectively, and

$$\theta = \frac{\lambda x_{d+1}}{\lambda x_{d+1} + (1 - \lambda)y_{d+1}} \in [0, 1].$$

Sweeping λ from 0 to 1 sweeps $\theta \in [0, 1]$, giving the result. \square

Based on Lemma B.1.9, we immediately obtain the following proposition.

Proposition B.1.10. *Let $C \subset \mathbb{R}^d \times \mathbb{R}_{++}$ be a convex set. Then $\text{pers}(C)$ is convex.*

Proof Let $x, y \in C$ and define $L = \{\lambda x + (1 - \lambda)y : \lambda \in [0, 1]\}$ to be the line segment between them. By Lemma B.1.9, $\text{pers}(L) = \{\lambda \text{pers}(x) + (1 - \lambda) \text{pers}(y) : \lambda \in [0, 1]\}$ is also a (convex) line segment, and we have $\text{pers}(L) \subset \text{pers}(C)$ as necessary. \square

B.1.2 Representation and separation of convex sets

JCD Comment: Put normal and tangent cones here

We now consider some properties of convex sets, showing that (1) they have nice separation properties—we can put hyperplanes between them—and (2) this allows several interesting representations of convex sets. We begin with the separation properties, developing them via the existence of projections. Interestingly, this existence of projections does not rely on any finite-dimensional structure, and can even be shown to hold in arbitrary Banach spaces (assuming the axiom of choice) [141]. We provide the results in a *Hilbert space*, meaning a complete vector space for which there exists an inner product $\langle \cdot, \cdot \rangle$ and associated norm $\|\cdot\|$ given by $\|x\|^2 = \langle x, x \rangle$. We first note that projections exist.

Theorem B.1.11 (Projections). *Let C be a closed convex set. Then for any x , there exists a unique point $\pi_C(x)$ minimizing $\|y - x\|$ over $y \in C$. Moreover, this point is characterized by the inequality*

$$\langle \pi_C(x) - x, y - \pi_C(x) \rangle \geq 0 \quad \text{for all } y \in C. \quad (\text{B.1.2})$$

Proof The existence and uniqueness of the projection follows from the parallelogram identity, that is, that for any x, y we have $\|x - y\|^2 + \|x + y\|^2 = 2(\|x\|^2 + \|y\|^2)$, which follows by noting that $\|x + y\|^2 = \|x\|^2 + \|y\|^2 + 2\langle x, y \rangle$. Indeed, let $\{y_n\} \subset C$ be a sequence such that

$$\|y_n - x\| \rightarrow \inf_{y \in C} \|y - x\| =: p_\star$$

as $n \rightarrow \infty$, where p_\star is the infimal value. We show that y_n is Cauchy, so that there exists a (unique) limit point of the sequence. Fix $\epsilon > 0$ and let N be such that $n \geq N$ implies $\|y_n - x\|^2 \leq p_\star^2 + \epsilon^2$. Let $m, n \geq N$. Then by the parallelogram identity,

$$\|y_n - y_m\|^2 = \|(x - y_n) - (x - y_m)\|^2 = 2 \left[\|x - y_n\|^2 + \|x - y_m\|^2 \right] - \|(x - y_n) + (x - y_m)\|^2.$$

Noting that

$$(x - y_n) + (x - y_m) = 2 \left[x - \frac{y_n + y_m}{2} \right] \quad \text{and} \quad \frac{y_n + y_m}{2} \in C \quad (\text{by convexity of } C),$$

we have

$$\|x - y_n\|^2 \leq p_\star^2 + \epsilon^2, \quad \|x - y_m\|^2 \leq p_\star^2 + \epsilon^2, \quad \text{and} \quad \|(x - y_n) + (x - y_m)\|^2 = 4 \left\| x - \frac{y_n + y_m}{2} \right\|^2 \geq 4p_\star^2.$$

In particular, we have

$$\|y_n - y_m\|^2 \leq 2 [p_\star^2 + \epsilon^2 + p_\star^2 + \epsilon^2] - 4p_\star^2 = 4\epsilon^2.$$

As $\epsilon > 0$ was arbitrary, this completes the proof of the first statement of the theorem.

To see the second result, assume that z is a point satisfying inequality (B.1.2), that is, such that

$$\langle z - x, y - z \rangle \geq 0 \quad \text{for all } y \in C.$$

Then we have

$$\|z - x\|^2 = \langle z - x, z - x \rangle = \underbrace{\langle z - x, z - y \rangle}_{\leq 0} + \langle z - x, y - x \rangle \leq \|z - x\| \|y - x\|$$

by the Cauchy-Schwarz inequality. Dividing both sides by $\|z - x\|$ yields $\|z - x\| \leq \|y - x\|$ for any $y \in C$, giving the result. Conversely, let $t \in [0, 1]$. Then for any $y \in C$,

$$\begin{aligned} \|\pi_C(x) - x\|^2 &\leq \|(1 - t)\pi_C(x) + ty - x\|^2 = \|\pi_C(x) - x + t(y - \pi_C(x))\|^2 \\ &= \|\pi_C(x) - x\|^2 + 2t\langle \pi_C(x) - x, y - \pi_C(x) \rangle + t^2 \|y - \pi_C(x)\|^2. \end{aligned}$$

Subtracting the projection value $\|\pi_C(x) - x\|^2$ from both sides and dividing by $t > 0$, we have

$$0 \leq 2\langle \pi_C(x) - x, y - \pi_C(x) \rangle + t \|y - \pi_C(x)\|^2.$$

Taking $t \rightarrow 0$ gives inequality (B.1.2). □

As an immediate consequence of Theorem B.1.11, we obtain several separation properties of convex sets, as well as a theorem stating that a closed convex set (not equal to the entire space in which it lies) can be represented as the intersection of all the half-spaces containing it.

Corollary B.1.12. *Let C be closed convex and $x \notin C$. Then there is a vector v strictly separating x from C , that is,*

$$\langle v, x \rangle > \sup_{y \in C} \langle v, y \rangle.$$

Moreover, we can take $v = x - \pi_C(x)$.

Proof By Theorem B.1.11, we know that taking $v = x - \pi_C(x)$ we have

$$0 \leq \langle y - \pi_C(x), \pi_C(x) - x \rangle = \langle y - \pi_C(x), -v \rangle = \langle y - x + v, -v \rangle = -\langle y, v \rangle + \langle x, v \rangle - \|v\|^2.$$

That is, we have $\langle v, y \rangle \leq \langle v, x \rangle - \|v\|^2$ for all $y \in C$ and $v \neq 0$. \square

Projections also never increase distances.

Corollary B.1.13. *Let C be closed convex and $y \in C$. Then for any x ,*

$$\|\pi_C(x) - y\| \leq \|x - y\|.$$

Proof Using inequality (B.1.2) in Theorem B.1.11, write

$$0 \geq \langle y - \pi_C(x), x - \pi_C(x) \rangle = \langle y - \pi_C(x), y - \pi_C(x) + x - y \rangle = \|y - \pi_C(x)\|^2 + \langle y - \pi_C(x), x - y \rangle.$$

Rearranging yields $\|y - \pi_C(x)\|^2 \leq \langle y - \pi_C(x), y - x \rangle \leq \|y - \pi_C(x)\| \|y - x\|$ by Cauchy-Schwarz. If $y = \pi_C(x)$, the result is trivial, and otherwise, dividing by $\|y - \pi_C(x)\| > 0$ gives the result. \square

In addition, we can show the existence of supporting hyperplanes, that is, hyperplanes “separating” the boundary of a convex set from itself.

Theorem B.1.14. *Let C be a convex set and $x \in \text{bd}(C)$, where $\text{bd}(C) = \text{cl}(C) \setminus \text{int } C$. Then there exists a non-zero vector v such that $\langle v, x \rangle \geq \sup_{y \in C} \langle v, y \rangle$.*

Proof Let $D = \text{cl}(C)$ be the closure of C and let $x_n \notin D$ be a sequence of points such that $x_n \rightarrow x$. Let us define the sequence of separating vectors $s_n = x_n - \pi_D(x_n)$ and the normalized version $v_n = s_n / \|s_n\|$. Notably, we have $\langle v_n, x_n \rangle > \sup_{y \in C} \langle v_n, y \rangle$ for all n . Now, the sequence $\{v_n\} \subset \{v : \|v\| = 1\}$ belongs to a compact set.¹ Passing to a subsequence if necessary, let us assume w.l.o.g. that $v_n \rightarrow v$ with $\|v\| = 1$. Then by a standard limiting argument for the $x_n \rightarrow x$, we have

$$\langle v, x \rangle \geq \langle v, y \rangle \text{ for all } y \in C,$$

which was our desired result. \square

JCD Comment: Picture of supporting hyperplanes and representations

Theorem B.1.14 gives us an important result. In particular, let D be an arbitrary set, and let $C = \text{cl Conv}(D)$ be the closure of the convex hull of D , which is the smallest closed convex set containing D . Then we can write C as the intersection of all the closed half-spaces containing D ; this is, in some sense, the most useful “convexification” of D . Recall that a closed half-space H is defined with respect to a vector v and real $a \in \mathbb{R}$ as

$$H := \{x : \langle v, x \rangle \leq a\}.$$

Before stating the theorem, we remark that by Observation B.1.6, the intersection of all the closed convex sets containing a set D is equal to the closure of the convex hull of D .

¹In infinite dimensions, this may not be the case. But we can apply the Banach-Alaoglu theorem, which states that, as v_n are linear operators, the sequence is weak-* compact, so that there is a vector v with $\|v\| \leq 1$ and a subsequence $m(n) \subset \mathbb{N}$ such that $\langle v_{m(n)}, x \rangle \rightarrow \langle v, x \rangle$ for all x .

Theorem B.1.15. *Let D be an arbitrary set. If $C = \text{cl Conv}(D)$, then*

$$C = \bigcap_{H \supset D} H, \quad (\text{B.1.3})$$

where H denotes a closed half-space containing D . Moreover, for any closed convex set C ,

$$C = \bigcap_{x \in \text{bd}(C)} H_x, \quad (\text{B.1.4})$$

where H_x denotes the intersection of halfspaces supporting C at x .

Proof We begin with the proof of the second result (B.1.4). Indeed, by Theorem B.1.14, we know that at each point x on the boundary of C , there exists a non-zero supporting hyperplane v , so that the half-space

$$H_{x,v} := \{y : \langle v, y \rangle \leq \langle v, x \rangle\} \supset C$$

is closed, convex, and contains C . We clearly have the containment $C \subset \bigcap_{x \in \text{bd}(C)} H_x$. Now let $x_0 \notin C$; we show that $x_0 \notin \bigcap_{x \in \text{bd}(C)} H_x$. As $x_0 \notin C$, the projection $\pi_C(x_0)$ of x_0 onto C satisfies $\langle x_0 - \pi_C(x_0), x_0 \rangle > \sup_{y \in C} \langle x_0 - \pi_C(x_0), y \rangle$ by Corollary B.1.12. Moreover, letting $v = x_0 - \pi_C(x_0)$, the hyperplane

$$h_{x_0,v} := \{y : \langle y, v \rangle = \langle \pi_C(x_0), v \rangle\}$$

is clearly supporting to C at the point $\pi_C(x_0)$. The half-space $\{y : \langle y, v \rangle \leq \langle \pi_C(x_0), v \rangle\}$ thus contains C and does not contain x_0 , implying that $x_0 \notin \bigcap_{x \in \text{bd}(C)} H_x$.

Now we show the first result (B.1.3). Let C be the closed convex hull of D and $T = \bigcap_{H \supset D} H$. By a trivial extension of the representation (B.1.4), we have that $C = \bigcap_{H \supset C} H$, where H denotes any halfspace containing C . As $C \supset D$, we have that $H \supset C$ implies $H \supset D$, so that

$$T = \bigcap_{H \supset D} H \subset \bigcap_{H \supset C} H = C.$$

On the other hand, as $C = \text{cl Conv}(D)$, Observation B.1.7 implies that any closed set containing D contains C . As a closed halfspace is convex and closed, we have that $H \supset D$ implies $H \supset C$, and thus $T = C$ as desired. \square

One elegant corollary of the closure operations and supporting hyperplanes for convex sets is that we can approximate convex hulls by expectations of vectors, even for potentially uncountable collections. By combining the strong law of large numbers with our descriptions of the convex hull, we have the following result.

Corollary B.1.16. *Let $X = \{x_\alpha\}_{\alpha \in \mathcal{A}} \subset \mathbb{R}^d$ be an arbitrary collection of vectors and let \mathcal{P} be the collection of probability distributions on elements of \mathcal{A} for which $\mathbb{E}_P[\|x_A\|] < \infty$, where $A \sim P$. Then*

$$\text{Conv}(X) = \{\mathbb{E}_P[x_A] \mid P \in \mathcal{P}\} \subset \text{cl Conv}(X).$$

If additionally X is compact (closed and bounded), then $\text{cl Conv}(X) = \{\mathbb{E}_P[x_A] \mid P \in \mathcal{P}\}$.

Proof Let $C = \{\mathbb{E}_P[x_A] \mid P \in \mathcal{P}\}$ be the middle set. We show that $\text{Conv}(X) \subset C \subset \text{cl Conv}(X)$. Taking any $P, Q \in \mathcal{P}$, we have $\lambda P + (1 - \lambda)Q \in \mathcal{P}$ for all $\lambda \in [0, 1]$, so that C is convex, giving $\text{Conv}(X) \subset C$. For the second inclusion, fix any $P \in \mathcal{P}$. Draw $A_1, A_2, \dots, A_n \stackrel{\text{iid}}{\sim} P$. Then by the strong law of large numbers, $\bar{x}_n := \frac{1}{n} \sum_{i=1}^n x_{A_i} \rightarrow \mathbb{E}_P[x_A]$ with probability 1, and so certainly there is a sequence of elements $\bar{x}_n \in \text{Conv}(X)$ satisfying $\bar{x}_n \rightarrow \mathbb{E}_P[x_A]$, and $\mathbb{E}_P[x_A] \in \text{cl Conv}(X)$.

To demonstrate that $\text{Conv}(X) = C$ requires more work. We prove the result by induction on the dimension. Consider the case that $d = 1$: then $\text{Conv}(X)$ takes on one of three forms: the open interval $\text{Conv}(X) = (b_0, b_1)$, the half-open interval $\text{Conv}(X) = (b_0, b_1]$ or $\text{Conv}(X) = [b_0, b_1)$, or the closed interval $\text{Conv}(X) = [b_0, b_1]$. In the last case, $\text{Conv}(X)$ is compact so that $\text{Conv}(X) = \text{cl Conv}(X)$ and we are done. Consider the first two, and w.l.o.g. assume $\mathbb{E}_P[x_A] = b_0$ while $x_\alpha > b_0$ for all $\alpha \in \mathcal{A}$. But for a distribution P on \mathcal{A} to yield $\mathbb{E}_P[x_A] = b_0$, it must be the case that $P(x_A = b_0) = 1$, a contradiction.

Now consider the d -dimensional case, and let $\mu = \mathbb{E}_P[x_A]$ for shorthand. Suppose for the sake of contradiction that that $\mu \in \text{cl Conv}(X) \setminus \text{Conv}(X)$. Then there is a non-zero vector v such that $\langle v, \mu \rangle \geq \langle v, \bar{x} \rangle$ for all $\bar{x} \in \text{Conv}(X)$. Letting $b = \langle v, \mu \rangle$, we have $\langle v, x_\alpha \rangle \leq b$ for all $\alpha \in \mathcal{A}$. That is, the hyperplane $H = \{x \mid \langle v, x \rangle = b\}$ separates μ from X , and the halfspace $H_- := \{x \mid \langle v, x \rangle \leq b\}$ contains X . So the scalar values $\langle v, x_\alpha \rangle \leq b$ for $\alpha \in \mathcal{A}$, and $\mathbb{E}_P[\langle v, x_A \rangle] = \langle v, \mu \rangle = b$ implies that $\langle v, x_A \rangle = b$ with probability 1 over $A \sim P$. In particular, the collection $\{x_\alpha \mid \langle v, x_\alpha \rangle = b\}$ is non-empty, so the sets $H \cap X$ and $H \cap \text{Conv}(X)$ are non-empty and of dimension $d - 1$. Induct downwards.

The final claim is simply Observation B.1.8. □

The second inclusion in Corollary B.1.16 can be strict even in one dimension: let $x_\alpha = \alpha$ for $\alpha \in (0, 1)$, so that $\text{Conv}(X) = X$, and any distribution P on α yields $\mathbb{E}_P[x_A] \in (0, 1)$.

B.2 Sublinear and support functions

A special case of convex functions will be sublinear functions, which form the basis of the transition between convex sets and convex functions. Accordingly, we give a special treatment here. Recall that f is convex if $f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$ for all x, y .

Definition B.6. A function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ is sublinear if it is convex and positively homogeneous, meaning

$$f(tx) = tf(x) \text{ for all } x \in \mathbb{R}^d \text{ and } t > 0.$$

Such functions are important in that they give some of the first dualities between convex sets and convex functions. As we see in the section to come, they also allow us to describe various first-order smoothness properties of convex functions.

The main result we shall need on sublinear functions is that they can be defined by a dual construction.

Proposition B.2.1. Let f be a closed sublinear function and define $S := \{s \mid \langle s, x \rangle \leq f(x) \text{ for all } x\}$. Then

$$f(x) = \sup_{s \in S} \langle s, x \rangle.$$

Proof As f is closed convex, there exist affine functions minorizing f at each point in its domain (Theorem B.3.3). That is, for some pair $(s, t) \in \mathbb{R}^d \times \mathbb{R}$, we have $\langle s, x \rangle - t \leq f(x)$ for all $x \in \mathbb{R}^d$.

Because necessarily $f(0) = 0$ by sublinearity, we have $t \geq 0$, and by positive homogeneity, we have $\langle s, \alpha x \rangle - t \leq f(\alpha x)$ for all $\alpha > 0$, that is, $\langle s, x \rangle - t/\alpha \leq f(x)$ for all x . Taking $\alpha \uparrow \infty$ we find that

$$\langle s, x \rangle \leq f(x) \text{ for all } x \in \mathbb{R}^d.$$

Because any closed convex function is the supremum of all affine functions minorizing it (Theorem B.3.7), we evidently have $f(x) = \sup_s \{\langle s, x \rangle \mid \langle s, \cdot \rangle \text{ minorizes } f\}$. \square

To any set S we can associate a particular sublinear function, the *support function* of S , defining

$$\sigma_S(x) := \sup_{s \in S} \langle s, x \rangle. \quad (\text{B.2.1})$$

This function is evidently a closed convex function—it is the supremum of linear functions—and is positively homogeneous, so that it is sublinear. We thus immediately have the duality

Corollary B.2.2. *Let f be a sublinear function. Then it is the support function of the closed convex set*

$$S_f := \{s \mid \langle s, x \rangle \leq f(x) \text{ for all } x \in \mathbb{R}^d\},$$

and hence if C is closed convex, then

$$C = \{x \mid \langle s, x \rangle \leq \sigma_C(s) \text{ for all } s \in \mathbb{R}^d\}.$$

A few other consequences of the definition are immediate. We see that σ_S has $\text{dom } \sigma_S = \mathbb{R}^d$ if and only if S is bounded: whenever $\|s\| \leq L$ for all $s \in S$, then $\sigma_S(x) \leq L\|x\|$. Conversely, if $\text{dom } \sigma_S = \mathbb{R}^d$ then it is locally Lipschitz (Theorem B.3.4) and (by positive homogeneity) thus globally Lipschitz, so we have $\langle s, x \rangle \leq \sigma_S(x) \leq L\|x\|$ for some $L < \infty$ and taking $x = s/\|s\|$ gives $\|s\| \leq L$. As another consequence, we see that support functions of a set S are the support functions of the closed convex hull of S :

Proposition B.2.3. *Let $S \subset \mathbb{R}^d$. Then*

$$\sigma_S(x) = \sigma_{\text{cl Conv } S}(x).$$

Proof Let $C = \text{Conv } S$, and let s_n be any sequence with $\langle s_n, x \rangle \rightarrow \sup_{s \in C} \langle s, x \rangle$. Then there exist $s_{n,i} \in S$, $i = 1, \dots, k(n)$, such that $s_n = \sum_{i=1}^{k(n)} \lambda_i s_{n,i}$ for some $\lambda \succeq 0$, $\langle \lambda, \mathbf{1} \rangle = 1$, which may change with n . But of course, $\langle s_n, x \rangle \leq \max_i \langle s_{n,i}, x \rangle$, and thus $\sigma_S(x) \geq \sigma_C(x)$. To see that $\sigma_C(x) = \sigma_{\text{cl } C}(x)$, note that for each $\epsilon > 0$, for each $s \in \text{cl } C$ there is $s' \in C$ with $\|s - s'\| < \epsilon$. Then $\langle s, x \rangle \leq \langle s', x \rangle + \epsilon\|x\|$ and $\sigma_{\text{cl } C}(x) \leq \sigma_C(x) + \epsilon\|x\|$. Take $\epsilon \downarrow 0$. \square

This proposition, coupled with Corollary B.2.2, shows that if sets S_1, S_2 have identical support functions, then they have identical closed convex hulls, and if they are closed convex, they are thus identical.

Corollary B.2.4. *Let $S_1, S_2 \subset \mathbb{R}^d$. If $\sigma_{S_1} = \sigma_{S_2}$, then $\text{cl Conv } S_1 = \text{cl Conv } S_2$.*

Proof By Proposition B.2.3, we have $\sigma_{S_i} = \sigma_{\text{cl Conv } S_i}$ for each i , and Corollary B.2.2 shows that if $\sigma_{C_1} = \sigma_{C_2}$ for closed convex sets C_1 and C_2 , then $C_1 = C_2$. \square

As another corollary, we have

Corollary B.2.5. *Let σ_1 and σ_2 be the support functions of the nonempty closed convex sets S_1 and S_2 . Then if $t_1 > 0$ and $t_2 > 0$,*

$$t_1\sigma_1 + t_2\sigma_2 = \sigma_{\text{cl}(t_1S_1+t_2S_2)}.$$

If either of S_1 or S_2 is compact, then $t_1\sigma_1 + t_2\sigma_2 = \sigma_{t_1S_1+t_2S_2}$.

Proof Let $S = t_1S_1 + t_2S_2$. In first statement, we have

$$\sigma_{\text{cl}S}(x) \stackrel{(\star)}{=} \sigma_S(x) = \sup \{ \langle t_1s_1 + t_2s_2, x \rangle \mid s_1 \in S_1, s_2 \in S_2 \},$$

equality (\star) following from Proposition B.2.3. As the suprema run independently through their respective sets S_1, S_2 , the latter quantity is evidently

$$\sigma_S(x) = t_1 \sup_{s_1 \in S_1} \langle s_1, x \rangle + t_2 \sup_{s_2 \in S_2} \langle s_2, x \rangle = t_1\sigma_{S_1}(x) + t_2\sigma_{S_2}(x).$$

The final result is an immediate consequence of the result that if C is a compact convex set and S is closed convex, then $C + S$ is closed convex. That $C + S$ is convex is immediate. To see that it is closed, let $x_n \in C, y_n \in S$ satisfy $x_n + y_n \rightarrow z$. Then proceeding to a subsequence, we have $x_{n(m)} \rightarrow x_\infty$ for some $x_\infty \in C$, and thus $y_{n(m)} \rightarrow z - x_\infty$, which is then necessarily in S . As the subsequence $x_{n(m)} + y_{n(m)} \rightarrow x_\infty + (z - x_\infty) \in C + S$ and $x_{n(m)} + y_{n(m)} \rightarrow z$ as well, this gives the result. \square

Linear transformations of support functions are also calculable. In the result, recall that for a matrix A and set S , the set $AS = \{As \mid s \in S\}$.

Proposition B.2.6. *Let $S \subset \mathbb{R}^d$ and $A \in \mathbb{R}^{m \times d}$. Then $\sigma_{\text{cl}AS}(x) = \sigma_S(A^\top x)$.*

Proof We have $\sigma_{AS}(x) = \sup_{s \in S} \langle As, x \rangle = \sup_{s \in S} \langle s, A^\top x \rangle$. The closure operation changes nothing (Proposition B.2.3). \square

Lastly, we show how to use support functions to characterize whether sets have interiors. Recall that for a set $S \subset \mathbb{R}^d$, the affine hull $\text{aff}(S)$ (Definition B.2) is the set of affine combinations of a point in S , and the relative interior of S is its interior relative to its affine hull (Definition B.4).

Proposition B.2.7. *Let $S \subset \mathbb{R}^d$ be non-empty a closed convex set. Then*

- (i) $s \in \text{int } S$ if and only if $\langle s, x \rangle < \sigma_S(x)$ for all $x \neq 0$.
- (ii) $s \in \text{relint } S$ if and only if $\langle s, x \rangle < \sigma_S(x)$ for all x with $\sigma_S(x) + \sigma_S(-x) > 0$.
- (iii) $\text{int } S$ is non-empty if and only if $\sigma_S(x) + \sigma_S(-x) > 0$ for all $x \neq 0$.

Proof

- (i) Because σ_S is positively homogeneous, an equivalent statement is that $\sigma_S(x) > \langle s, x \rangle$ for all $x \in \mathbb{S}^{d-1} = \{x \in \mathbb{R}^d \mid \|x\|_2 = 1\}$. If $s \in \text{int } S$, we there exists $\epsilon > 0$ such that $s + \epsilon x \in S$ for all $x \in \mathbb{S}^{d-1}$, and so

$$\sigma_S(x) \geq \langle s + \epsilon x, x \rangle = \langle s, x \rangle + \epsilon,$$

so that $\langle s, x \rangle < \sigma_S(x)$.

Conversely, let s be any point satisfying $\sigma_S(x) - \langle s, x \rangle > 0$ for all $x \in \mathbb{S}^{d-1}$. Because σ_S is lower semicontinuous, the infimum $\inf_{x \in \mathbb{S}^{d-1}} \{\sigma_S(x) - \langle s, x \rangle\}$ is attained at some $x^* \in \mathbb{S}^{d-1}$ (see Proposition C.0.1). Then there exists some $\epsilon > 0$ such that $\langle s, x \rangle + \epsilon \leq \sigma_S(x)$ for all $x \in \mathbb{S}^{d-1}$. Let u be any vector with $\|u\|_2 < \epsilon$. Then $\langle s + u, x \rangle = \langle s, x \rangle + \langle u, x \rangle \leq \langle s, x \rangle + \epsilon \leq \sigma_S(x)$, so Corollary B.2.2 implies $s + u \in S$ and $s \in \text{int } S$.

- (ii) We decompose \mathbb{R}^d into subspaces $V \oplus U$, where $U = V^\perp$ and V is parallel to $\text{aff}(S)$. Writing $x = x_U + x_V$, where $x_U \in U$ and $x_V \in V$, the function $\langle s, x_U \rangle$ is constant for $s \in S$. Repeat the argument for part (i) in the subspace V .
- (iii) Suppose $\text{int } S$ is non-empty. Then $s \in \text{int } S$ implies $\langle s, x \rangle < \sigma_S(x)$ for all x with $\|x\| = 1$. Then $\sigma_S(x) + \sigma_S(-x) > \langle s, x - x \rangle = 0$. Conversely, if $\text{int } S$ is empty, there exists a hyperplane containing S (by a dimension counting argument and that the relative interior of S is never empty [111, Theorem III.2.1.3]), which we may write as $S \subset \{s \mid v^T s = b\}$ for some $v \neq 0$. For this $\sigma_S(v) + \sigma_S(-v) = b - b = 0$.

□

B.3 Convex functions

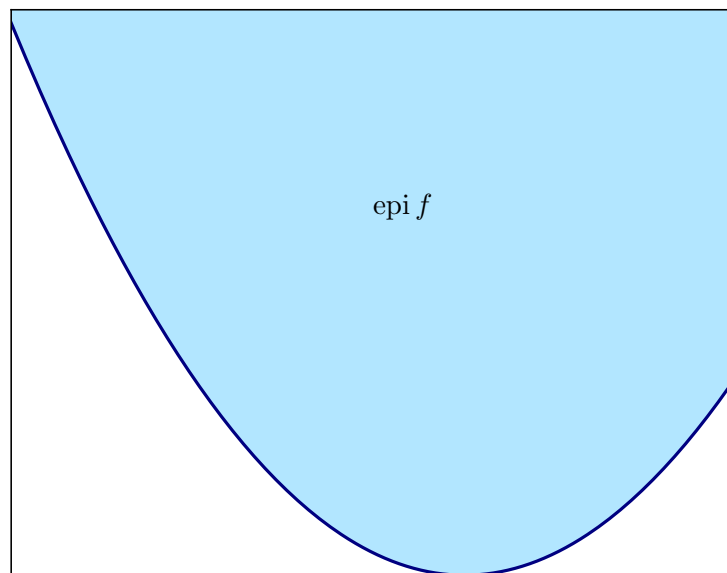


Figure B.1: The epigraph of a convex function.

We now build off of the definitions of convex sets to define convex functions. As we will see, convex functions have several nice properties that follow from the geometric (separation) properties of convex sets. First, we have

Definition B.7. A function f is convex if for all $\lambda \in [0, 1]$ and $x, y \in \text{dom } f$,

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y). \quad (\text{B.3.1})$$

We define the domain $\text{dom } f$ of a convex function to be those points x such that $f(x) < +\infty$. Note that Definition B.7 implies that the domain of f must be convex.

An equivalent definition of convexity follows by considering a natural convex set attached to the function f , known as its epigraph.

Definition B.8. The epigraph $\text{epi } f$ of a function is the set

$$\text{epi } f := \{(x, t) : t \in \mathbb{R}, f(x) \leq t\}.$$

That is, the epigraph of a function f is the set of points on or above the graph of the function itself, as depicted in Figure B.1. It is immediate from the definition of the epigraph that f is convex if and only if $\text{epi } f$ is convex. Thus, we see that any convex set $C \subset \mathbb{R}^{d+1}$ that is unbounded “above,” meaning that $C = C + \{0\} \times \mathbb{R}_+$, defines a convex function, and conversely, any convex function defines such a set C . This duality in the relationship between a convex function and its epigraph is central to many of the properties we exploit.

B.3.1 Equivalent definitions of convex functions

We begin our discussion of convex functions by enumerating a few standard properties that also characterize convexity. The simplest of these relate to properties of the derivatives and second derivatives of functions. We begin by elucidating one of the most basic properties of convexity: that the slopes of convex functions are increasing. Beginning with functions on \mathbb{R} , suppose that $f : \mathbb{R} \rightarrow \mathbb{R}$ is convex, and let $x \in \text{dom } f$ and $v \in \mathbb{R}$ be otherwise arbitrary. Then define the quotient function

$$q(t) := \frac{f(x + tv) - f(x)}{t}, \quad t \geq 0, \quad (\text{B.3.2})$$

which we claim is nondecreasing in $t \geq 0$ if and only if f is convex. Indeed, let $t \geq s > 0$ and define $\lambda = \frac{s}{t} \in [0, 1]$. Then

$$\begin{aligned} q(t) \geq q(s) & \text{ if and only if } \lambda[f(x + tv) - f(x)] \geq f(x + \lambda tv) - f(x) \\ & \text{ if and only if } \lambda f(x + tv) + (1 - \lambda)f(x) \geq f((1 - \lambda)x + \lambda(x + tv)), \end{aligned}$$

the latter holding for all λ if and only if f is convex.

JCD Comment: Draw a picture of increasing quotient

Because the quotient function (B.3.2) is nondecreasing, we can relatively straightforwardly give first-order characterizations of convexity as well. Indeed, suppose that $f : \mathbb{R} \rightarrow \mathbb{R}$ is differentiable; then convexity is equivalent to the first-order inequality that for all $x, y \in \mathbb{R}$, we have

$$f(y) \geq f(x) + f'(x)(y - x). \quad (\text{B.3.3})$$

To see that inequality (B.3.3) implies that f is convex follows from algebraic manipulations: let $\lambda \in [0, 1]$ and $z = \lambda x + (1 - \lambda)y$, so that $y - z = \lambda(y - x)$ and $x - z = (1 - \lambda)(x - y)$. Then

$$f(y) \geq f(z) + \lambda f'(z)(y - x) \quad \text{and} \quad f(x) \geq f(z) + (1 - \lambda)f'(z)(x - y),$$

and multiplying the former by $(1 - \lambda)$ and the latter by λ and adding the two inequalities yields

$$\lambda f(x) + (1 - \lambda)f(y) \geq \lambda f(z) + (1 - \lambda)f(z) + \lambda(1 - \lambda)f'(z)(y - x) + \lambda(1 - \lambda)f'(z)(x - y) = f(\lambda x + (1 - \lambda)y),$$

as desired. Conversely, let $v = y - x$ in the quotient (B.3.2), so that $q(t) = \frac{f(x+tv) - f(x)}{t}$, which is non-decreasing. If f is differentiable, we see that $q(0) := \lim_{t \downarrow 0} q(t) = f'(x)(y - x)$, and so

$$q(1) = f(y) - f(x) \geq q(0) = f'(x)(y - x)$$

as desired.

We may also give the standard second order characterization: if $f : \mathbb{R} \rightarrow \mathbb{R}$ is twice differentiable and $f''(x) \geq 0$ for all x , then f is convex. To see this, note that

$$f(y) = f(x) + f'(x)(y - x) + \frac{1}{2}f''(tx + (1 - t)y)(x - y)^2$$

for some $t \in [0, 1]$ by Taylor's theorem, so that $f(y) \geq f(x) + f'(x)(y - x)$ for all x, y because $f''(tx + (1 - t)y) \geq 0$. As a consequence, we obtain inequality (B.3.3), which implies that f is convex.

As convexity is a property that depends only on properties of functions on lines—one dimensional projections—we can straightforwardly extend the preceding results to functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$. Indeed, noting that if $h(t) = f(x + ty)$ then $h'(0) = \langle \nabla f(x), y \rangle$ and $h''(0) = y^\top \nabla^2 f(x) y$, we have that a differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) \quad \text{for all } x, y,$$

while a twice differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if

$$\nabla^2 f(x) \succeq 0 \quad \text{for all } x.$$

Noting that nothing in the derivation that the quotient (B.3.2) was non-decreasing relied on f being a function on \mathbb{R} , we can see that a function $f : \mathbb{R}^d$ is convex if and only if it satisfies the *increasing slopes* criterion: for all $x \in \text{dom } f$ and any vector v , the quotient

$$t \mapsto q(t) := \frac{f(x + tv) - f(x)}{t} \tag{B.3.4}$$

is nondecreasing in $t \geq 0$ (where we leave x, v implicit). An alternative version of the criterion (B.3.4) is that if $x \in \text{dom } f$ and v is any vector, if we define the one-dimensional convex function $h(t) = f(x + tv)$ then for any $s < t$ and $\Delta > 0$, we have

$$\frac{h(t + \Delta) - h(t)}{\Delta} \geq \frac{h(t) - h(s)}{t - s} \geq \frac{h(t) - h(s - \Delta)}{t - (s - \Delta)}. \tag{B.3.5}$$

The proof that either of the inequalities (B.3.5) is equivalent to convexity we leave as an exercise (Q. C.1).

JCD Comment: Draw pictures of increasing slopes

We summarize each of these implications in a theorem for reference.

Proposition B.3.1 (Convexity). *The following are all equivalent:*

(i) The function f is convex.

(ii) The function f satisfies the criterion of increasing slopes (B.3.4).

If f is differentiable (respectively, twice differentiable), the following are also equivalent:

(iii) The function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle \quad \text{for all } x, y.$$

(iv) The function f has positive semidefinite Hessian: $\nabla^2 f(x) \succeq 0$ for all x .

JCD Comment: Draw a picture and of strict convexity

A condition slightly stronger than convexity is *strict convexity*, which makes each of the inequalities in Proposition B.3.1 strict. We begin with the classical definition: a function f is strictly convex if it is convex and

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y)$$

whenever $\lambda \in (0, 1)$ and $x \neq y \in \text{dom } f$. These are convex functions, but always have strictly increasing slopes—secants lie strictly above f . By tracing through the arguments leading to Proposition B.3.1 (replace appropriate non-strict inequalities with strict inequalities), one obtains the following corollary describing strictly convex functions.

Corollary B.3.2 (Strict convexity). *The following are all equivalent:*

(i) The function f is strictly convex.

(ii) The function f has strictly increasing slopes (B.3.4).

If f is differentiable (respectively, twice differentiable), the following are also equivalent:

(iii) The function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies

$$f(y) > f(x) + \langle \nabla f(x), y - x \rangle \quad \text{for all } x \neq y.$$

(iv) The function f has positive definite Hessian: $\nabla^2 f(x) \succ 0$ for all x .

B.3.2 Continuity properties of convex functions

We now consider a few continuity properties of convex functions and a few basic relationships of the function f to its epigraph. First, we give a definition of the *subgradient* of a convex function.

Definition B.9. A vector g is a subgradient of f at a point x_0 if for all x ,

$$f(x) \geq f(x_0) + \langle g, x - x_0 \rangle. \tag{B.3.6}$$

The subdifferential or subgradient set of f at x_0 is

$$\partial f(x_0) := \{g \mid f(x) \geq f(x_0) + \langle g, x - x_0 \rangle \text{ for all } x\}.$$

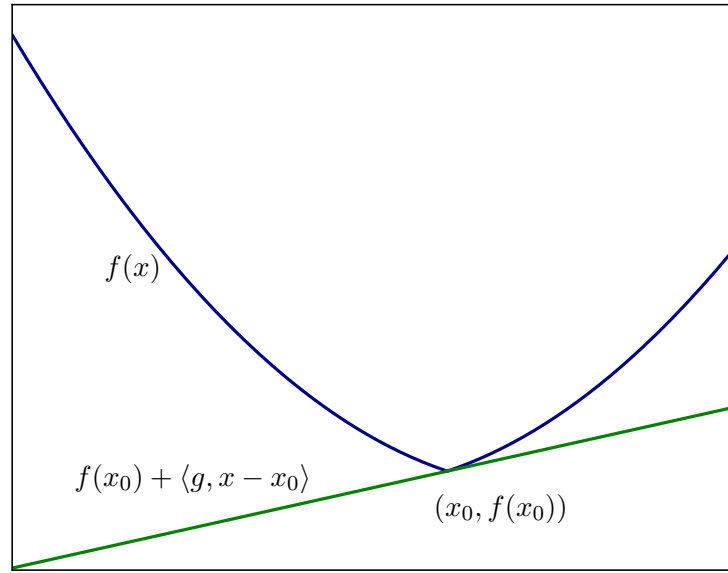


Figure B.2. The tangent (affine) function to the function f generated by a subgradient g at the point x_0 .

See Figure B.2 for an illustration of the affine minorizing function given by the subgradient of a convex function at a particular point.

Interestingly, convex functions have subgradients (at least, nearly everywhere). This is perhaps intuitively obvious by viewing a function in conjunction with its epigraph $\text{epi } f$ and noting that $\text{epi } f$ has supporting hyperplanes, but here we state a result that will have further use.

Theorem B.3.3. *Let f be convex. Then there is an affine function minorizing f . More precisely, for any $x_0 \in \text{relint dom } f$, there exists a vector g such that*

$$f(x) \geq f(x_0) + \langle g, x - x_0 \rangle.$$

Proof If $\text{relint dom } f = \emptyset$, then it is clear that f is either identically $+\infty$ or its domain is a single point $\{x_0\}$, in which case the constant function $f(x_0)$ minorizes f . Now, we assume that $\text{int dom } f \neq \emptyset$, as we can simply always change basis to work in the affine hull of $\text{dom } f$.

We use Theorem B.1.14 on the existence of supporting hyperplanes to construct a subgradient. Indeed, we note that $(x_0, f(x_0)) \in \text{bd epi } f$, as for any open set O we have that $(x_0, f(x_0)) + O$ contains points both inside and outside of $\text{epi } f$. Thus, Theorem B.1.14 guarantees the existence of a vector v and $a \in \mathbb{R}$, not both simultaneously zero, such that

$$\langle v, x_0 \rangle + af(x_0) \leq \langle v, x \rangle + at \quad \text{for all } (x, t) \in \text{epi } f. \quad (\text{B.3.7})$$

Inequality (B.3.7) implies that $a \geq 0$, as for any x we may take $t \rightarrow +\infty$ while satisfying $(x, t) \in \text{epi } f$. Now we argue that $a > 0$ strictly. To see this, note that for suitably small $\delta > 0$, we have $x = x_0 - \delta v \in \text{dom } f$. Then we find by inequality (B.3.7) that

$$\langle v, x_0 \rangle + af(x_0) \leq \langle v, x_0 \rangle - \delta \|v\|^2 + af(x_0 - \delta v), \quad \text{or} \quad a[f(x_0) - f(x_0 - \delta v)] \leq -\delta \|v\|^2.$$

So if $v = 0$, then Theorem B.1.14 already guarantees $a \neq 0$, while if $v \neq 0$, then $\|v\|^2 > 0$ and we must have $a \neq 0$ and $f(x_0) \neq f(x_0 - \delta v)$. As we showed already that $a \geq 0$, we must have $a > 0$.

Then by setting $t = f(x_0)$ and dividing both sides of inequality (B.3.7) by a , we obtain

$$\frac{1}{a} \langle v, x_0 - x \rangle + f(x_0) \leq f(x) \quad \text{for all } x \in \text{dom } f.$$

Setting $g = -v/a$ gives the result of the theorem, as we have $f(x) = +\infty$ for $x \notin \text{dom } f$. \square

Convex functions generally have quite nice behavior. Indeed, they enjoy some quite remarkable continuity properties just by virtue of the defining convexity inequality (B.3.1). In particular, the following theorem shows that convex functions are continuous on the relative interiors of their domains. Even more, convex functions are Lipschitz continuous on any compact subsets contained in the (relative) interior of their domains. (See Figure B.3 for an illustration of this fact.)

Theorem B.3.4. *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and $C \subset \text{relint dom } f$ be compact. Then there exists an $L = L(C) \geq 0$ such that*

$$|f(x) - f(x')| \leq L \|x - x'\|.$$

As an immediate consequence of Theorem B.3.4, we note that if $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex and defined everywhere on \mathbb{R}^d , then it is continuous. Moreover, we also have that $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is continuous everywhere on the (relative) interior of its domain: let any $x_0 \in \text{relint dom } f$. Then for small enough $\epsilon > 0$, the set $\text{cl}(\{x_0 + \epsilon B\} \cap \text{dom } f)$, where $B = \{x : \|x\|_2 \leq 1\}$, is a closed and bounded—and hence compact—set contained in the (relative) interior of $\text{dom } f$. Thus f is Lipschitz on this set, which is a neighborhood of x_0 . In addition, if $f : \mathbb{R} \rightarrow \mathbb{R}$, then f is continuous everywhere except (possibly) at the endpoints of its domain.

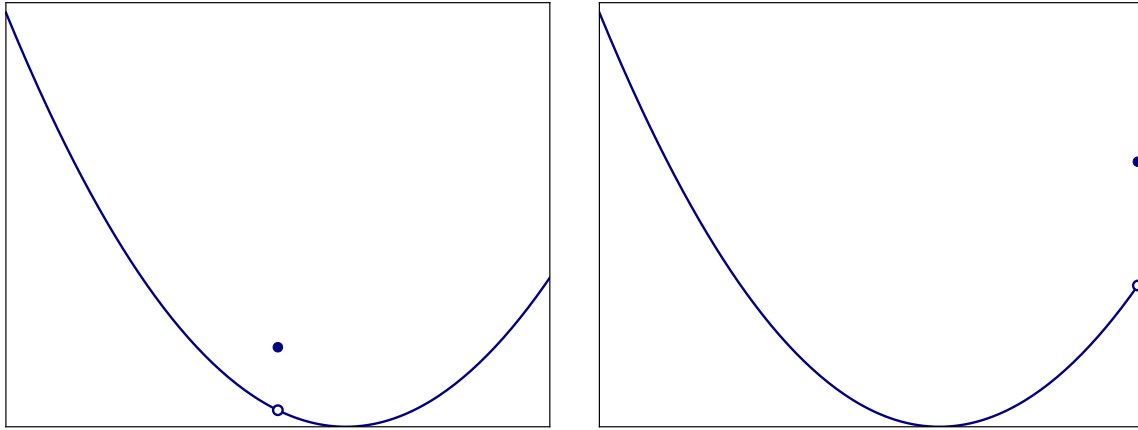


Figure B.3. Left: discontinuities in $\text{int dom } f$ are impossible while maintaining convexity (Theorem B.3.4). Right: At the edge of $\text{dom } f$, there may be points of discontinuity.

Proof of Theorem B.3.4 To prove the theorem, we require a technical lemma.

Lemma B.3.5. *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and suppose that there are x_0 , $\delta > 0$, m , and M such that*

$$m \leq f(x) \leq M \quad \text{for } x \in B(x_0, 2\delta) := \{x : \|x - x_0\| < 2\delta\}.$$

Then f is Lipschitz on $B(x_0, \delta)$, and moreover,

$$|f(y) - f(y')| \leq \frac{M - m}{\delta} \|y - y'\| \quad \text{for } y, y' \in B(x_0, \delta).$$

Proof Let $y, y' \in B(x_0, \delta)$, and define $y'' = y' + \delta(y' - y)/\|y' - y\| \in B(x_0, 2\delta)$. Then we can write y' as a convex combination of y and y'' , specifically,

$$y' = \frac{\|y' - y\|}{\delta + \|y' - y\|} y'' + \frac{\delta}{\delta + \|y' - y\|} y.$$

Thus we obtain by convexity

$$\begin{aligned} f(y') - f(y) &\leq \frac{\|y' - y\|}{\delta + \|y' - y\|} f(y'') + \frac{\delta}{\delta + \|y' - y\|} f(y) - f(y) = \frac{\|y - y'\|}{\delta + \|y - y'\|} [f(y'') - f(y)] \\ &\leq \frac{M - m}{\delta + \|y - y'\|} \|y - y'\|. \end{aligned}$$

Here we have used the bounds on f assumed in the lemma. Swapping the assignments of y and y' gives the same lower bound, thus giving the desired Lipschitz continuity. \square

With Lemma B.3.5 in place, we proceed to the proof proper. We assume without loss of generality that $\text{dom } f$ has an interior; otherwise we prove the theorem restricting ourselves to the affine hull of $\text{dom } f$. The proof follows a standard compactification argument. Suppose that for each $x \in C$, we could construct an open ball $B_x = B(x, \delta_x)$ with $\delta_x > 0$ such that

$$|f(y) - f(y')| \leq L_x \|y - y'\| \quad \text{for } y, y' \in B_x. \quad (\text{B.3.8})$$

As the B_x cover the compact set C , we can extract a finite number of them, call them B_{x_1}, \dots, B_{x_k} , covering C , and then within each (overlapping) ball f is $\max_k L_{x_k}$ Lipschitz. As a consequence, we find that

$$|f(y) - f(y')| \leq \max_k L_{x_k} \|y - y'\|$$

for any $y, y' \in C$.

We thus must derive inequality (B.3.8), for which we use the boundedness Lemma B.3.5. We must demonstrate that f is bounded in a neighborhood of each $x \in C$. To that end, fix $x \in \text{int dom } f$, and let the points x_0, \dots, x_d be affinely independent and such that

$$\Delta := \text{Conv}\{x_0, \dots, x_d\} \subset \text{dom } f$$

and $x \in \text{int } \Delta$; let $\delta > 0$ be such that $B(x, 2\delta) \subset \Delta$. Then by Carathéodory's theorem (Theorem B.1.3) we may write any point $y \in B(x, 2\delta)$ as $y = \sum_{i=0}^d \lambda_i x_i$ for $\sum_i \lambda_i = 1$ and $\lambda_i \geq 0$, and thus

$$f(y) \leq \sum_{i=0}^d \lambda_i f(x_i) \leq \max_{i \in \{0, \dots, d\}} f(x_i) =: M.$$

Moreover, Theorem B.3.3 implies that there is some affine h function minorizing f ; let $h(x) = a + \langle v, x \rangle$ denote this function. Then

$$m := \inf_{x \in C} f(x) \geq \inf_{x \in C} h(x) = a + \inf_{x \in C} \langle v, x \rangle > -\infty$$

exists and is finite, so that in the ball $B(x, 2\delta)$ constructed above, we have $f(y) \in [m, M]$ as required by Lemma B.3.5. This guarantees the existence of a ball B_x required by inequality (B.3.8). \square

Our final discussion of continuity properties of convex functions revolves around the most common and analytically convenient type of convex function, the so-called *closed-convex* functions.

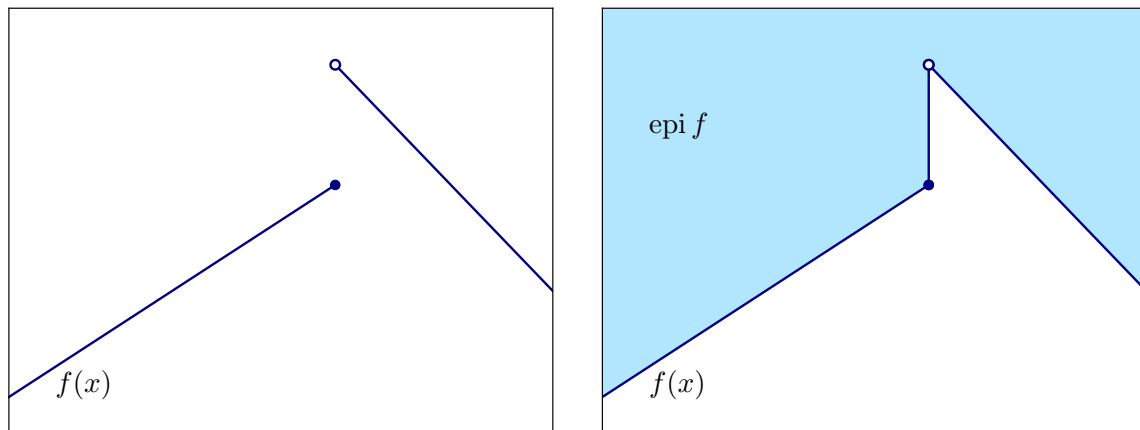


Figure B.4. A closed—equivalently, lower semi-continuous—function. On the right is shown the *closed* epigraph of the function.

Definition B.10. A function f is closed if its epigraph, $\text{epi } f$, is a closed set.

Equivalently, a function is closed if it is lower semi-continuous, meaning that

$$\liminf_{x \rightarrow x_0} f(x) \geq f(x_0) \quad (\text{B.3.9})$$

for all x_0 and any sequence of points tending toward x_0 . See Figure B.4 for an example such function and its associated epigraph.

Interestingly, in the one-dimensional case, closed convexity implies continuity. Indeed, we have the following observation (compare Figures B.4 and B.3 previously):

Observation B.3.6. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a closed convex function. Then f is continuous on its domain, and for any $x_0 \in \text{bd dom } f$, $\lim_{x \rightarrow x_0} f(x) = f(x_0)$ whether or not $x_0 \in \text{dom } f$.

Proof By Theorem B.3.4, we need only consider the endpoints of the domain of f (the result is obvious by Theorem B.3.4 if $\text{dom } f = \mathbb{R}$); let $x_0 \in \text{bd dom } f$. Let $y \in \text{dom } f$ be an otherwise arbitrary point, and define $x = \lambda y + (1 - \lambda)x_0$. Then taking $\lambda \rightarrow 0$, we have

$$f(x) \leq \lambda f(y) + (1 - \lambda)f(x_0) \rightarrow f(x_0),$$

so that $\limsup_{x \rightarrow x_0} f(x) \leq f(x_0)$. By the closedness assumption (B.3.9), we have $\liminf_{x \rightarrow x_0} f(x) \geq f(x_0)$, and continuity follows. Note that in this argument, if $x_0 \notin \text{dom } f$, then $f(x_0) = +\infty$ by convention; for $\text{epi } f$ to be closed we require that for each $t < f(x_0) = \infty$, we may take a small enough open interval $U = (y, x_0)$ for which $f(x) > t$ for all $x \in U$. \square

In the full-dimensional case, we do not have quite the same continuity, though Theorem B.3.4 guarantees continuity on the (relative) interior of $\text{dom } f$.

An important characterization of convex functions is as the supremum of all affine functionals (linear plus an offset) below them, which is one of the keys to duality relationships about functions to come.

Theorem B.3.7. Let f be closed convex and let \mathcal{A} be the collection of affine functions h satisfying $f(x) \geq h(x)$ for all x . Then $f(x) = \sup_{h \in \mathcal{A}} h(x)$.

Proof By Theorem B.1.15 that any closed convex set is the intersection of all the halfspaces containing (even supporting) it, we can write $\text{epi } f = \bigcap_{H \in \mathcal{H}} H$, where \mathcal{H} is the collection of closed halfspaces $H \supset \text{epi } f$. We may write any such halfspace as

$$H = \{(x, r) \in \mathbb{R}^d \times \mathbb{R} \mid \langle a, x \rangle + br \leq c\}$$

where $(a, b) \in \mathbb{R}^d \times \mathbb{R}$ is non-zero. As $H \supset \text{epi } f$, the particular nature of epigraphs (that is, that if $(x, t) \in \text{epi } f$ then $(x, t + \Delta) \in \text{epi } f$ for all $\Delta > 0$) means that $b \leq 0$, and so for any $b < 0$ we may divide through by b to rewrite H as $H = \{(x, r) \mid \langle a/b, x \rangle + r \geq c/b\}$, while if $b = 0$ then $H = \{(x, r) \mid \langle a, x \rangle \leq c\}$. That is, it is no loss of generality to set

$$\begin{aligned}\mathcal{H}_1 &:= \{\text{Halfspaces } \{(x, r) \mid \langle a, x \rangle + r \geq c\} \text{ containing } \text{epi } f\} \\ \mathcal{H}_0 &:= \{\text{Halfspaces } \{(x, r) \mid \langle a, x \rangle \geq c\} \text{ containing } \text{epi } f\},\end{aligned}$$

which (respectively) correspond to the non-vertical halfspaces containing $\text{epi } f$ and the halfspaces containing $\text{dom } f \subset \mathbb{R}^d$. We have $\text{epi } f = \bigcap_{H \in \mathcal{H}_1} H \cap \bigcap_{H \in \mathcal{H}_0} H$.

Identify the halfspaces $H \in \mathcal{H}_0$ or \mathcal{H}_1 with the associated triple $(a, 0, c)$ or $(a, 1, c)$ and abuse notation to write $(a, i, c) \in \mathcal{H}_i$ for $i \in \{0, 1\}$. For any $(a, 1, c) \in \mathcal{H}_1$, the linear function

$$l(x) = c - \langle a, x \rangle = \inf\{r \mid \langle a, x \rangle + r \geq c\} \text{ satisfies } \langle a, x \rangle + l(x) \geq c \text{ for all } x,$$

and so necessarily $l(x) \leq f(x)$ for all x , while for the function $h(x) = \sup_{(a, 1, c) \in \mathcal{H}_1} \{c - \langle a, x \rangle\}$ we have

$$\text{epi } h = \bigcap_{H \in \mathcal{H}_1} H.$$

Thus, if we can show that

$$\bigcap_{H \in \mathcal{H}_1} H \cap \bigcap_{H \in \mathcal{H}_0} H = \bigcap_{H \in \mathcal{H}_1} H \tag{B.3.10}$$

the proof will be complete.

To show the equality (B.3.10), take arbitrary vectors $v_0 = (a_0, 0, c_0) \in \mathcal{H}_0$ and $v_1 = (a_1, 1, c_1) \in \mathcal{H}_1$, and let $H_0 = \{(x, r) \mid \langle a_0, x \rangle \geq c_0\}$ and $H_1 = \{(x, r) \mid \langle a_1, x \rangle + r \geq c_1\}$ be the associated halfspaces. Consider the conic-like vector

$$v(t) := (a_1 + ta_0, 1, c_1 + tc_0) \quad \text{for } t \geq 0$$

and associated halfspace $H(t) := \{(x, r) \mid \langle a_1 + ta_0, x \rangle + r \geq c_1 + tc_0\}$. Then as $\langle a_0, x \rangle \geq c_0$ if and only if $t\langle a_0, x \rangle \geq tc_0$ for all $t \geq 0$, any point $(x, r) \in H_0 \cap H_1$ satisfies

$$\langle a_1 + ta_0, x \rangle + r \geq c_1 + tc_0 \quad \text{for } t \geq 0,$$

that is, $H(t) \in \mathcal{H}_1$ and $(x, r) \in \bigcap_{t \geq 0} H(t)$. Additionally, taking $t = 0$ we see that $H(0) = H_1$ and so $\bigcap_{t \geq 0} H(t) \subset H_1$, while taking $t \uparrow \infty$ we obtain that each $(x, r) \in \bigcap_{t \geq 0} H(t)$ satisfies $\langle a_0, x \rangle \geq c_0$. That is, we have

$$\bigcap_{t \geq 0} H(t) = H_0 \cap H_1,$$

while $H(t) \in \mathcal{H}_1$ for all $t \geq 0$. This shows the equality (B.3.10). \square

JCD Comment: Show a picture of the above argument

In spite of the continuity of closed convex functions on \mathbb{R} , closed convex functions on higher dimensional spaces need not be continuous. Indeed, it is immediate (see Proposition B.3.9 to follow) that $f(x) := \sup_{\alpha \in \mathcal{A}} \{f_\alpha(x)\}$ is closed convex whenever f_α are all closed convex for any index set \mathcal{A} . We have the following failure of continuity.

Example B.3.8 (A discontinuous closed convex function): Define the function $f : \mathbb{R}^2 \rightarrow \overline{\mathbb{R}}$ by

$$f(x) := \sup \left\{ \alpha x_1 + \beta x_2 \mid \frac{1}{2} \alpha^2 \leq \beta \right\}.$$

Then certainly $f(\mathbf{0}) = 0$ and f is closed convex. If the supremum is attained then $\beta = \frac{1}{2} \alpha^2$ and so $\beta \geq 0$ and

$$f(x) = \sup_{\alpha} \left\{ \alpha x_1 + \frac{1}{2} \alpha^2 x_2 \right\} = \begin{cases} 0 & \text{if } x = \mathbf{0} \\ -\frac{x_1^2}{2x_2} & \text{if } x_2 < 0 \\ +\infty & \text{otherwise.} \end{cases}$$

But then along the path $x_2 = -\frac{1}{2} x_1^2$, we always have $f(x) = 1$, while taking $x_1 \rightarrow 0$ gives $f(x) = 1 > 0 = f(\mathbf{0})$. \diamond

B.3.3 Operations preserving convexity

We now turn to a description of a few simple operations on functions that preserve convexity. First, we extend the intersection properties of convex sets to operations on convex functions. (See Figure B.5 for an illustration of the proposition.)

Proposition B.3.9. *Let $\{f_\alpha\}_{\alpha \in \mathcal{A}}$ be an arbitrary collection of convex functions indexed by \mathcal{A} . Then*

$$f(x) := \sup_{\alpha \in \mathcal{A}} f_\alpha(x)$$

is convex. Moreover, if for each $\alpha \in \mathcal{A}$, the function f_α is closed convex, f is closed convex.

Proof The proof is immediate once we consider the epigraph $\text{epi } f$. We have that

$$\text{epi } f = \bigcap_{\alpha \in \mathcal{A}} \text{epi } f_\alpha,$$

which is convex whenever $\text{epi } f_\alpha$ is convex for all α and closed whenever $\text{epi } f_\alpha$ is closed for all α (recall Observation B.1.5). \square

Another immediate result is that composition of a convex function with an affine transformation preserves convexity:

Proposition B.3.10. *Let $A \in \mathbb{R}^{d \times n}$ and $b \in \mathbb{R}^d$, and let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex. Then the function $g(y) = f(Ay + b)$ is convex.*

Partial minimization of convex functions and some related transformations preserve convexity as well.

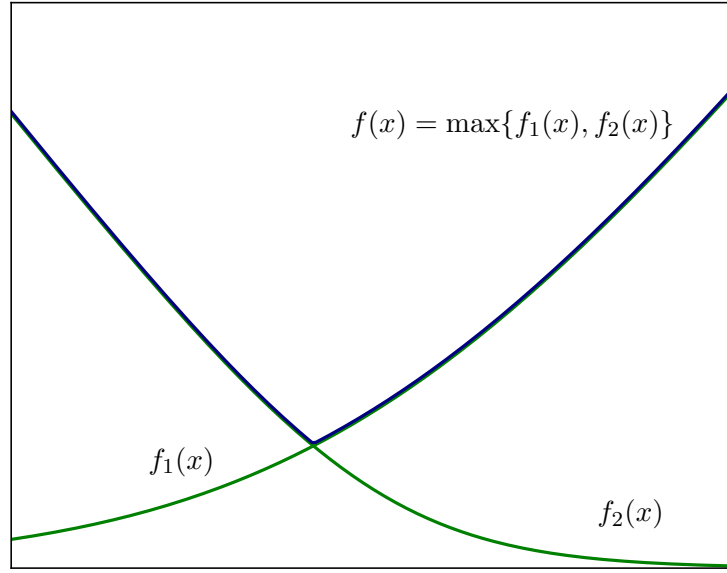


Figure B.5. The maximum of two convex functions is convex, as its epigraph is the intersection of the two epigraphs.

Proposition B.3.11. Let $A \in \mathbb{R}^{d \times n}$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex, and $Y \subset \mathbb{R}^d$ be convex. Then $g(x) = \inf\{f(y) \mid Ay = x, y \in Y\}$ is convex. If Y is compact and f is closed convex, then g is closed convex.

Proof Let $x_0, x_1 \in \mathbb{R}^n$. If $Ay = x_0$ has no solution in $y \in Y$, then $g(x_0) = +\infty$, and similarly if $Ay = x_1$ has no solutions then $g(x_1) = +\infty$, and we trivially have $g(\lambda x_0 + (1 - \lambda)x_1) \leq +\infty$ in either case for all $\lambda \in (0, 1)$. Assuming that the sets $\{y \in Y \mid Ay = x_0\}$ and $\{y \in Y \mid Ay = x_1\}$ are non-empty, let $\epsilon > 0$ be arbitrary and y_0, y_1 satisfy $Ay_i = x_i$ and that $f(y_i) \leq g(x_i) + \epsilon$. Then $y_\lambda = \lambda y_0 + (1 - \lambda)y_1$ satisfies $Ay_\lambda = \lambda x_0 + (1 - \lambda)x_1$, and so

$$g(\lambda x_0 + (1 - \lambda)x_1) \leq f(\lambda y_0 + (1 - \lambda)y_1) \leq \lambda f(y_0) + (1 - \lambda)f(y_1) \leq \lambda g(x_0) + (1 - \lambda)g(x_1) + \epsilon$$

for all $\lambda \in [0, 1]$. Take $\epsilon \rightarrow 0$.

For the lower semicontinuity (closed convexity) statement, let $x_n \rightarrow x$; we wish to show that $\liminf_n g(x_n) \geq g(x)$. If $g(x_n) = +\infty$ for all x_n , then we trivially have the result. Otherwise, assume $g(x_n) < \infty$ for all n , let $\epsilon > 0$ be arbitrary, and let $y_n \in Y$ satisfy $Ay_n = x_n$ and $f(y_n) \leq g(x_n) + \epsilon$. Then as Y is compact, y_n has convergent subsequences; let y be any such limit. We have $Ay = x$, and $g(x) \leq f(y) \leq \liminf_n f(y_n) \leq \liminf_n g(x_n) + \epsilon$. As $\epsilon > 0$ was arbitrary, we have the result. \square

From the proposition we immediately see that if $f(x, y)$ is jointly convex in x and y , then the partially minimized function $\inf_{y \in Y} f(x, y)$ is convex whenever Y is a convex set.

Lastly, we consider the functional analogue of the perspective transform. Given a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, the *perspective transform* of f is defined as

$$\text{pers}(f)(x, t) := \begin{cases} tf\left(\frac{x}{t}\right) & \text{if } t > 0 \text{ and } \frac{x}{t} \in \text{dom } f \\ +\infty & \text{otherwise.} \end{cases} \quad (\text{B.3.11})$$

In analogue with the perspective transform of a convex set, the perspective transform of a function is (jointly) convex.

Proposition B.3.12. *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex. Then $\text{pers}(f) : \mathbb{R}^{d+1} \rightarrow \mathbb{R}$ is convex.*

Proof The result follows if we can show that $\text{epi pers}(f)$ is a convex set. With that in mind, note that

$$\mathbb{R}^d \times \mathbb{R}_{++} \times \mathbb{R} \ni (x, t, r) \in \text{epi pers}(f) \text{ if and only if } f\left(\frac{x}{t}\right) \leq \frac{r}{t}.$$

Rewriting this, we have

$$\begin{aligned} \text{epi pers}(f) &= \left\{ (x, t, r) \in \mathbb{R}^d \times \mathbb{R}_{++} \times \mathbb{R} : f\left(\frac{x}{t}\right) \leq \frac{r}{t} \right\} \\ &= \left\{ t(x', 1, r') : x' \in \mathbb{R}^d, t \in \mathbb{R}_{++}, r' \in \mathbb{R}, f(x') \leq r' \right\} \\ &= \{t(x, 1, r) : t > 0, (x, r) \in \text{epi } f\} = \mathbb{R}_{++} \times \{(x, 1, r) : (x, r) \in \text{epi } f\}. \end{aligned}$$

This is a convex cone. □

Finally, we discuss closing a convex function, that is, replacing f with $\text{cl } f$, the lower semicontinuous closure of f . For $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$, we define $\text{cl } f$ pointwise by

$$\text{cl } f(x) := \liminf_{y \rightarrow x} f(y) \tag{B.3.12}$$

to be the lower semicontinuous closure of f . (Recall Definition B.10.) We demonstrate that this function is indeed the largest closed function below f .

Lemma B.3.13. *The function $\text{cl } f$ is lower-semicontinuous, $\text{epi cl } f = \text{cl epi } f$, and if g is any closed function with $g \leq f$, then $g \leq \text{cl } f$.*

Proof Assuming the second statement, the first follows trivially, so let us prove the second. Let $(x, r) \in \text{cl epi } f$. Then there exists a sequence $(x_n, r_n) \in \text{epi } f$ such that $(x_n, r_n) \rightarrow (x, r)$, and $f(x_n) \leq r_n$. Then evidently $\text{cl } f(x) \leq \liminf r_n$, so that $(x, r) \in \text{epi cl } f$. Conversely, let $(x, r) \in \text{epi cl } f$. Let $x_n \rightarrow x$ satisfy $f(x_n) \rightarrow \liminf_{y \rightarrow x} f(y) = \text{cl } f(x)$. Then because $r \geq \text{cl } f(x)$, there exist $r_n \geq f(x_n)$ for which $r_n \rightarrow r$, while $(x_n, r_n) \in \text{epi } f$. In particular, $(x_n, r_n) \rightarrow (x, r)$, which is thus in $\text{cl epi } f$.

Finally, we prove the third claim of the lemma. If g is closed, then $\text{epi } g$ is closed as well, and $\text{epi } g \supset \text{epi } f$. Then $\text{epi } g \supset \text{cl epi } f = \text{epi cl } f$, that is, $\text{cl } f(x) \leq r$ implies $g(x) \leq r$, i.e., $\text{cl } f \geq g$. □

Using this lemma, we can provide some additional color to Theorem B.3.7 that a closed convex function is equal to the supremum of the affine functions minorizing it. In particular, Observation B.1.6 that closures of convex sets remain convex implies that if f is convex, then $\text{cl } f$ is convex as well. We collect this and a bit more with the following proposition, which now immediately follows.

Proposition B.3.14. *Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be a convex function. Let \mathcal{A} be the collection of all affine functions $h \leq f$. Then $\text{cl } f(x) = \sup_{h \in \mathcal{A}} h(x)$ for all x . In particular, if f is lower semicontinuous at x , then $f(x) = \sup_{h \in \mathcal{A}} h(x)$.*

B.3.4 Smoothness properties, first-order developments for convex functions, and subdifferentiability

In addition to their continuity properties, convex functions typically enjoy strong differentiability properties. Some of these interact with the duality properties we present in the section C.2 to follow. Our main goal will be to show how there exist (roughly) derivative-like objects for convex functions, so that for some suitably nice object $D_f(x, v)$ we have

$$f(x + tv) = f(x) + D_f(x, v)t + o(t) \quad (\text{B.3.13})$$

for t small and any v . In the case that f is differentiable, of course, this must coincide with the usual derivative, so that $D_f(x, v) = \langle \nabla f(x), v \rangle$. For convex functions, a directional derivative *always* exists (even if f is non-differentiable), meaning that we can make sense of the first-order development (B.3.13) in some generality.

As one prototypical result, we leverage Rademacher's theorem on almost everywhere differentiability of Lipschitz functions to show that convex functions are almost everywhere differentiable:

Theorem B.3.15 (Rademacher). *Let $U \subset \mathbb{R}^d$ be open and $f : U \rightarrow \mathbb{R}^k$ be Lipschitz continuous. Then f is differentiable almost everywhere on U .*

Proofs of this result are standard in measure-theoretic analysis texts; see, e.g., [89, Section 3.5] or [188, Theorem 10.8(ii)]. As any convex function is locally Lipschitz on its domain (recall Theorem B.3.4), we thus have the following result (where we assume that $\text{dom } f$ has an interior).

Corollary B.3.16. *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex. Then it is differentiable except on a set of Lebesgue measure zero on its domain.*

Other differentiability properties of convex functions are also of interest. We begin by considering directional differentiability properties, after which we expand to consider differentiability and continuous differentiability of (convex) functions. To begin, recall that the *directional derivative* of a function f in direction v at x is

$$f'(x; v) := \lim_{t \downarrow 0} \frac{f(x + tv) - f(x)}{t} \quad (\text{B.3.14})$$

when this quantity exists. When $f'(x; v)$ exists for all directions v and is linear in v , we call the function *Gateaux differentiable*. A (stronger in infinite dimensions) notion of differentiability is *Fréchet differentiability*: f has Fréchet differential g at x if

$$f(y) = f(x) + \langle g, y - x \rangle + o(\|y - x\|) \quad (\text{B.3.15})$$

as $y \rightarrow x$, which is then uniform in the distance $\|y - x\|$. It is immediate that if f is Fréchet differentiable with derivative g then it is Gateaux differentiable with $f'(x; v) = \langle g, v \rangle$. Conveniently, in finite dimensions, these notions coincide with the standard gradient, and $f'(x; v) = \langle \nabla f(x), v \rangle$, whenever f is locally Lipschitzian.

Proposition B.3.17. *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be Gateaux differentiable at x , that is, its directional derivative $f'(x; v)$ is linear in v , and locally Lipschitz, so that there exists $L < \infty$ such that $|f(x) - f(y)| \leq L\|x - y\|$ for y near x . Then f is Fréchet differentiable with Fréchet differential*

$$\nabla f(x) = \left[\frac{\partial f(x)}{\partial x_j} \right]_{j=1}^d,$$

and $f'(x; v) = \langle \nabla f(x), v \rangle$ and $\|\nabla f(x)\| \leq L$.

Proof If f is Fréchet differentiable at x with differential g , then we immediately have

$$\frac{f(x + tv) - f(x)}{t} = \frac{t\langle g, v \rangle + o(t)}{t} \rightarrow \langle g, v \rangle$$

as $t \rightarrow 0$, so that it is Gateaux differentiable.

Conversely, suppose that $f'(x; v) = \langle g, v \rangle$ for all $v \in \mathbb{R}^d$ for some $g \in \mathbb{R}^d$. Assume for the sake of contradiction that f is not Fréchet differentiable at x , so that

$$\limsup_{\|\Delta\| \downarrow 0} \frac{f(x + \Delta) - f(x) - \langle g, \Delta \rangle}{\|\Delta\|} = c > 0.$$

Take any sequence $\Delta_n \rightarrow 0$ achieving this limit supremum, and let $\Delta_n = \epsilon_n v_n$ for a sequence v_n on the sphere, that is, $\|v_n\| = 1$, so $\epsilon_n = \|\Delta_n\|$. Then by passing to a subsequence if necessary, we can assume w.l.o.g. that $v_n \rightarrow v$ with $\|v\| = 1$. Then

$$\begin{aligned} \frac{|f(x + \Delta_n) - f(x) - \langle g, \Delta_n \rangle|}{\epsilon_n} &= \frac{|f(x + \epsilon_n v + \epsilon_n(v_n - v)) - f(x) - \epsilon_n \langle g, v \rangle - \epsilon_n \langle g, v_n - v \rangle|}{\epsilon_n} \\ &\leq \frac{|f(x + \epsilon_n v) - f(x) - \epsilon_n \langle g, v \rangle|}{\epsilon_n} + \frac{L \epsilon_n \|v_n - v\| + \epsilon_n \|g\| \|v_n - v\|}{\epsilon_n}. \end{aligned}$$

Both of these terms tend to zero, a contradiction, and so f is Fréchet differentiable at x , and its Fréchet derivative is g . That Fréchet differentiability implies differentiability follows by noting that the partial derivatives $f'(x; e_j) = \frac{\partial f(x)}{\partial x_j}$ for each coordinate j .

Finally, the Lipschitzian bound on $\|\nabla f(x)\|$ follows by noting that

$$L \|\Delta\| \geq |f(x + \Delta) - f(x)| = |\langle \nabla f(x), \Delta \rangle| + o(\|\Delta\|).$$

Taking $\Delta = tv$ and $t \downarrow 0$, this implies that $L \|v\| \geq \langle \nabla f(x), v \rangle$ for all v , which is equivalent to $\|\nabla f(x)\| \leq L$. \square

The main consequence of convexity that is important for us is that a convex function is directionally differentiable at every point in the interior of its domain, though the directional derivative need not be linear:

Proposition B.3.18. *Let f be convex and $x \in \text{int dom } f$. Then $f'(x; v)$ exists and the mapping $v \mapsto f'(x; v)$ is sublinear, convex, and globally Lipschitz.*

Proof If $x \in \text{int dom } f$, then the criterion (B.3.4) of increasing slopes guarantees that $f'(x; v) = \lim_{t \downarrow 0} \frac{f(x + tv) - f(x)}{t}$ exists for all $x \in \text{int dom } f$, as the quantity is monotone. To see that $f'(x; v)$ is convex and sublinear in v , note that positive homogeneity is immediate, as we have $\frac{1}{t}(f(x + \alpha tv) - f(x)) = \frac{\alpha}{\alpha t}(f(x + \alpha tv) - f(x))$ for all $\alpha > 0$, and $f'(x; 0) = 0$. That it is convex is straightforward as well: for any u, v we have

$$\frac{f(x + t(\lambda u + (1 - \lambda)v)) - f(x)}{t} \leq \lambda \frac{f(x + tu) - f(x)}{t} + (1 - \lambda) \frac{f(x + tv) - f(x)}{t}$$

and take $t \downarrow 0$. For the global Lipschitz claim, note that f is already locally Lipschitz near $x \in \text{int dom } f$ (recall Theorem B.3.4), so that there exists some $L < \infty$ and $\epsilon > 0$ such that for all $\|v\| = 1$ and $0 \leq t \leq \epsilon$ $|f(x + tv) - f(x)| \leq Lt$, whence $|f'(x; v)| \leq L$ and by homogeneity

$$|f'(x; v)| \leq L \|v\| \text{ for all } v. \quad \square$$

An inspection of the proof shows that the result extends even to all of $\text{dom } f$ if we allow $f'(x; v) = +\infty$ whenever $x + tv \notin \text{dom } f$ for all $t > 0$, though of course we lose that $f'(x; v)$ is finite-valued. Then we have the following corollary, showing that $f'(x; v)$ provides a valid first-order development of f in all directions from x (where we take $\infty \cdot t = \infty$ whenever $t > 0$).

Corollary B.3.19. *Let $x \in \text{dom } f$. Then*

$$f(x + tv) = f(x) + f'(x; v)t + o(t)$$

as $t \downarrow 0$ and

$$f(x + tv) \geq f(x) + f'(x; v)t \text{ for all } t \geq 0.$$

Proof The first part is immediate by definition of $f'(x; v) = \lim_{t \downarrow 0} \frac{f(x+tv) - f(x)}{t}$. The second is immediate from the criterion (B.3.4) of increasing slopes, as the limit in the directional derivative (B.3.14) becomes an infimum for convex functions: $f'(x; v) = \inf_{t > 0} \frac{f(x+tv) - f(x)}{t}$. \square

There are strong connections between subdifferentials and directional derivatives, and hence of the local developments (B.3.13). The following result makes this clear.

Proposition B.3.20. *Let f be convex and $x \in \text{relint dom } f$. Then*

$$\partial f(x) = \{s \mid \langle s, v \rangle \leq f'(x; v) \text{ for all } v\} \neq \emptyset.$$

Proof For shorthand let $S = \{s \mid \langle s, v \rangle \leq f'(x; v) \text{ all } v\}$ be the set on the right. If $s \in S$, then the criterion (B.3.4) of increasing slopes guarantees that

$$\langle s, v \rangle \leq \frac{f(x + tv) - f(x)}{t} \text{ for all } v \in \mathbb{R}^d, t > 0.$$

Recognizing that as v is allowed to vary over all of \mathbb{R}^d and $t > 0$, then $x + tv$ similarly describes \mathbb{R}^d , we see that this condition is completely equivalent to the definition (B.3.6) of the subgradient.

That $\partial f(x) \neq \emptyset$ is Theorem B.3.3. \square

We can also extend this to $x \in \text{dom } f$ —not necessarily the interior—where we see that there is no loss (even when f may be $+\infty$ valued) to defining

$$\partial f(x) := \{s \mid \langle s, v \rangle \leq f'(x; v) \text{ for all } v\}. \quad (\text{B.3.16})$$

Notably, the directional derivative function $v \mapsto f'(x; v)$ always exists for $x \in \text{dom } f$ and is a sublinear convex function, and thus $\partial f(x)$ above is always a closed convex set whose support function (recall (B.2.1)) is the closure of $v \mapsto f'(x; v)$. While the subdifferential $\partial f(x)$ is always a compact convex set when $x \in \text{int dom } f$, even when it exists it may not be compact if x is on the boundary of $\text{dom } f$. To see one important example of this, consider the indicator function

$$\mathbf{I}_C(x) := \begin{cases} +\infty & \text{if } x \notin C \\ 0 & \text{if } x \in C \end{cases}$$

of a closed convex set C . For simplicity, let $C = [a, b]$ be an interval. Then we have

$$\partial \mathbf{I}_C(x) = \begin{cases} [0, \infty] & \text{if } x = b \\ \{0\} & \text{if } a < x < b \\ [-\infty, 0] & \text{if } x = a. \end{cases}$$

Whether points $\pm\infty$ are included is a matter of convenience and whether we work with the extended real line.

JCD Comment: Draw a picture of this

These representations points to a certain closure property of subgradients, namely, that the subdifferential is closed under additions of the normal cone to the domain of f :

Lemma B.3.21. *Let $\mathcal{N}_{\text{dom } f}(x)$ be the normal cone (Definition C.1) to $\text{dom } f$ at the point x (where $\mathcal{N}_{\text{dom } f}(x) = \{0\}$ for $x \in \text{int dom } f$ and $\mathcal{N}_{\text{dom } f}(x) = \emptyset$ for $x \notin \text{dom } f$). Then*

$$\partial f(x) = \partial f(x) + \mathcal{N}_{\text{dom } f}(x).$$

In particular, if x is a boundary point $x \in \text{bd dom } f$ of the domain of f , then either $\partial f(x) = \emptyset$ or $\partial f(x)$ is unbounded.

Proof We only need concern ourselves with points $x \in \text{bd dom } f$, where the normal cone $\mathcal{N} = \mathcal{N}_{\text{dom } f}(x)$ is non-trivial. If $\partial f(x)$ is empty, there is nothing to prove, so assume that $\partial f(x)$ is non-empty. Then the definition (B.3.16) of the subdifferential as $\partial f(x) = \{s \mid \langle s, u \rangle \leq f'(x; u)\}$ allows us to prove the result. First, consider vectors u for which $f'(x; u) = +\infty$. Then certainly, for any $s \in \partial f(x)$, we have $\langle s + v, u \rangle \leq f'(x; u)$ for all $v \in \mathcal{N}$. If $f'(x; u) < \infty$, then for small enough $t > 0$ we necessarily have $x + tu \in \text{dom } f$. In particular, the definition of the normal cone gives that $v \in \mathcal{N}$ satisfies $0 \geq \langle v, x + tu - x \rangle = t \langle v, u \rangle$, or that $\langle v, u \rangle \leq 0$. Thus $\langle s + v, u \rangle \leq \langle s, u \rangle \leq f'(x; u)$, and so $s + v \in \partial f(x)$ once again.

The claim about boundedness is immediate, because $\mathcal{N}_{\text{dom } f}$ is a cone. □

A more compelling case for the importance of the subgradient set with respect to first-order developments and differentiability properties of convex functions is the following:

JCD Comment: Add a picture of this as well.

Proposition B.3.22. *Let f be convex and $x \in \text{int dom } f$. Then*

$$\begin{aligned} f(y) &= f(x) + \sup_{s \in \partial f(x)} \langle s, y - x \rangle + o(\|y - x\|) \\ &= f(x) + f'(x; y - x) + o(\|y - x\|). \end{aligned}$$

Proof That $\sup_{s \in \partial f(x)} \langle s, v \rangle = f'(x; v)$ is immediate by Theorem B.3.7 and Proposition B.2.1, because $f'(x; v)$ is sublinear and closed convex in v when $x \in \text{int dom } f$. Certainly the right hand sides are then equal.

We thus prove the equality $f(y) = f(x) + f'(x; y - x) + o(\|y - x\|)$, where the argument is similar to that for Proposition B.3.17. Let $y_n \rightarrow x$ be any sequence and let $\Delta_n = y_n - x$, so that $\|\Delta_n\| \rightarrow 0$; as $x \in \text{int dom } f$, there exists a (local) Lipschitz constant L such that $|f(x + \Delta) - f(x)| \leq L \|\Delta\|$

for all small Δ . Similarly, because $v \mapsto f'(x; v)$ is convex (even positively homogeneous and thus sublinear), it has a Lipschitz constant, and we take this to be L as well. Now, write $\Delta_n = \epsilon_n v_n$ where $\|v_n\| = 1$ and $\epsilon_n \rightarrow 0$, and moving to a subsequence if necessary let $v_n \rightarrow v$. Then we have

$$\begin{aligned} f(x + \Delta_n) - f(x) - f'(x; \Delta_n) &= f(x + \epsilon_n v + \epsilon_n(v_n - v)) - f(x) - f'(x; \epsilon_n v + \epsilon_n(v_n - v)) \\ &= f(x + \epsilon_n v) - f(x) - f'(x; \epsilon_n v) \pm 2L\epsilon_n \|v_n - v\| \\ &= o(\epsilon_n) \end{aligned}$$

because $\|v_n - v\| \rightarrow 0$ and $f(x + \epsilon_n v) - f(x) = f'(x; \epsilon_n v) + o(\epsilon_n)$ by definition of the directional derivative. \square

Note that convexity only played the role of establishing the local Lipschitz property of f in the proof of Proposition B.3.22; any locally Lipschitz function with directional derivatives will enjoy a similar first-order expansion.

As our final result on smoothness properties of convex functions, we connect subdifferentials to differentiability properties of convex f . First, we give a lemma showing that the subdifferential set ∂f is outer semicontinuous.

Lemma B.3.23 (Closure of the graph of the subdifferential). *Let $f : \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ be closed convex. Then the graph $\{(x, s) \mid x \in \mathbb{R}^d, s \in \partial f(x)\}$ of its subdifferential is closed. Equivalently, whenever $x_n \rightarrow x$ with $s_n \in \partial f(x_n)$ and $s_n \rightarrow s$, f has non-empty subdifferential at x with $s \in \partial f(x)$.*

Proof We prove the second statement, whose equivalence to the first is definitional. Fix any $y \in \mathbb{R}^d$. Then $f(y) \geq f(x_n) + \langle s_n, y - x_n \rangle$, and because f is closed (i.e., lower semicontinuous), we have $\liminf f(x_n) \geq f(x)$. Let $\epsilon > 0$ be arbitrary. Then for all large enough n , we have $f(x_n) \geq f(x) - \epsilon$, and similarly, $\|s_n - s\| \leq \epsilon$, $\|x_n - x\| \leq \epsilon$, and $\|y - x_n\| \leq \|y - x\| + \epsilon$. Then

$$\begin{aligned} f(y) &\geq f(x_n) + \langle s_n, y - x_n \rangle \geq f(x) + \langle s, y - x \rangle - \epsilon - \|s - s_n\| \|y - x_n\| \\ &\geq f(x) + \langle s, y - x \rangle - \epsilon - \epsilon \|y - x_n\| - \|s\| \|x - x_n\| \\ &\geq f(x) + \langle s, y - x \rangle - \epsilon - \epsilon(1 + \epsilon) \|y - x\| - \|s\| \epsilon. \end{aligned}$$

As ϵ was arbitrary we have $f(y) \geq f(x) + \langle s, y - x \rangle$ as desired. \square

Given the somewhat technical Lemma B.3.23, we can show that if f is convex and differentiable at a point, it is in fact continuously differentiable at the point.

Proposition B.3.24. *Let f be convex and $x \in \text{int dom } f$. Then $\partial f(x)$ is a singleton if and only if f is differentiable at x . If additionally f is differentiable on an open set U , then f is continuously differentiable on U .*

Proof Because $x \in \text{int dom } f$, there exists $L < \infty$ such that f is L -Lipschitz near x by Theorem B.3.4. Suppose that $\partial f(x) = \{s\}$. Then the directional derivative $f'(x; v) = \langle s, v \rangle$ for all v , and Proposition B.3.22 gives

$$f(y) = f(x) + \langle s, y - x \rangle + o(\|y - x\|)$$

as $y \rightarrow x$, that is, f is differentiable. Conversely, assume that f is differentiable at x . Then taking any vector v , we immediately have $f'(x; v) = \langle \nabla f(x), v \rangle$ and Proposition B.3.20 gives that $\partial f(x) = \{\nabla f(x)\}$.

To see that f is in fact continuously differentiable on U , let $x \in U$ and f be L -Lipschitz on a compact set $C \subset U$ containing x in its interior. Let $x_k \in C$ satisfy $x_k \rightarrow x$ and let $s_k = \nabla f(x_k) \in \partial f(x_k)$. Then $\|s_k\| \leq L$, and each subsequence has a further convergent subsequence. Lemma B.3.23 implies that any convergent subsequence $s_{k(m)} \rightarrow s \in \partial f(x)$. But as $\partial f(x) = \{\nabla f(x)\}$, we have $\nabla f(x_{k(m)}) \rightarrow \nabla f(x)$ and so $\nabla f(x)$ is continuous in x . \square

B.3.5 Calculus rules of subgradients

We close this section with a few calculus results on subdifferentials of convex functions. These calculus rules show that the subdifferential plays a similar role to the gradient for differentiable functions. Additionally, they allow us to take derivatives of various extremal functions.

Our first result shows that subdifferentials of sums are sums of subdifferentials, which relies on both the representation of sublinear functions as support functions for convex sets and the characterization of the subdifferential in terms of directional derivatives:

Proposition B.3.25. *Let f and g be closed convex functions, and let $x \in \text{int dom } f$ and g be subdifferentiable at x , meaning that $\partial g(x) \neq \emptyset$. Then*

$$\partial(f + g)(x) = \partial f(x) + \partial g(x).$$

Proof By Proposition B.3.20, the set $\partial f(x)$ is a compact convex set, and the general definition (B.3.16) of the subdifferential gives that $\partial g(x)$ is closed convex. Let $S_1 = \partial f(x)$ and $S_2 = \partial g(x)$. Then immediately $S_1 + S_2 \subset \partial(f + g)(x)$, so that

$$S := \partial(f + g)(x) = \left\{ s \mid \langle s, v \rangle \leq f'(x; v) + g'(x; v) \text{ for all } v \in \mathbb{R}^d \right\}$$

is non-empty. Because of the support function equality $f'(x; v) = \sigma_{S_1}(v)$ and $g'(x; v) = \sigma_{S_2}(v)$, Corollary B.2.5 gives

$$\sigma_S(v) = \sigma_{S_1}(v) + \sigma_{S_2}(v) = \sigma_{S_1 + S_2}(v).$$

Thus (Corollary B.2.4) $S_1 + S_2 = S$. \square

Other situations that arise frequently are composition with affine mappings and taking maxima or suprema of convex functions, so that finding a calculus for these is also important.

Corollary B.3.26. *Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$ be convex and for $A \in \mathbb{R}^{m \times d}$ and $b \in \mathbb{R}^m$, let $g(x) = f(Ax + b)$. Then*

$$\partial g(x) = A^T \partial f(Ax + b).$$

Proof Using the directional derivative, we have $g'(x; v) = f'(Ax + b; Av)$ for all $v \in \mathbb{R}^d$, and applying Proposition B.2.6 gives that the latter is the support function of the convex compact set $A^T \partial f(Ax + b)$. \square

It is also useful to be able to compute subdifferentials of maxima and suprema (recall Proposition B.3.9). Consider a collection $\{f_\alpha\}_{\alpha \in \mathcal{A}}$ of convex functions, and define

$$f(x) := \sup_{\alpha \in \mathcal{A}} f_\alpha(x). \tag{B.3.17}$$

The function f is certainly convex. For a given x let

$$\mathcal{A}(x) := \{\alpha \in \mathcal{A} \mid f_\alpha(x) = f(x)\}$$

be the indices attaining the suprema, that is, the active index set (though this may be empty). Then there is an “easy” direction:

Lemma B.3.27. *With the notation above,*

$$\partial f(x) \supset \text{cl Conv} \left\{ \bigcup \partial f_\alpha(x) \mid \alpha \in \mathcal{A}(x) \right\} = \text{cl Conv} \{g \mid g \in \partial f_\alpha(x) \text{ for some } \alpha \in \mathcal{A}(x)\}.$$

Proof Let $\alpha \in \mathcal{A}(x)$ and $g \in \partial f_\alpha(x)$. Then

$$f(y) \geq f_\alpha(y) \geq f_\alpha(x) + \langle g, y - x \rangle = f(x) + \langle g, y - x \rangle.$$

Thus $g \in \partial f(x)$, which as a closed convex set must thus include its closed convex hull. \square

A much more challenging argument is to show that the active index set $\mathcal{A}(x)$ exactly characterizes the subdifferential of f at x ; we simply state a typical result as a proposition.

Proposition B.3.28. *Let \mathcal{A} be a compact set (for some metric) and assume that for each x , the mapping $\alpha \mapsto f_\alpha(x)$ is upper semi-continuous. Then*

$$\partial f(x) = \text{Conv} \left\{ \bigcup \partial f_\alpha(x) \mid \alpha \in \mathcal{A}(x) \right\} = \text{Conv} \{g \mid g \in \partial f_\alpha(x) \text{ for some } \alpha \in \mathcal{A}(x)\}.$$

For a proof, see [111, Theorem 4.4.2].

JCD Comment: Draw a picture of this

Finally, we revisit the partial minimization operation in Proposition B.3.11. In this case, we require a bit more care when defining subdifferentials and subdifferentiability. For $A \in \mathbb{R}^{n \times m}$ with $m \geq n$, where A has rank n (so that $x \mapsto Ax$ is surjective) and $f : \mathbb{R}^m \rightarrow \bar{\mathbb{R}}$, define the function

$$f_A(x) = \inf \{f(y) \mid Ay = x\},$$

which is convex. Define the set $Y^*(x) := \{y \mid Ay = x \text{ and } f_A(x) = f(y)\}$ to be the set of y attaining the infimum in the definition of f_A , which may be empty. When it is not, however, we can characterize the subdifferential of $f_A(x)$:

Proposition B.3.29. *Let $x \in \mathbb{R}^n$ be a point for which $Y^*(x)$ is non-empty for the function f_A . Then*

$$\partial f_A(x) = \{s \mid A^T s \in \partial f(y)\}$$

for any $y \in Y^*(x)$, and the set on the right is independent of the choice of y .

Proof A vector s is a subgradient of f at x if and only if

$$f_A(x') \geq f_A(x) + \langle s, x' - x \rangle \text{ for all } x' \in \mathbb{R}^n,$$

which (as $Ay = x$ for $y \in Y^*(x)$) is equivalent to

$$f_A(x') \geq f(y) + \langle s, x' - Ay \rangle \text{ for all } x' \in \mathbb{R}^n.$$

Because A has full row rank, for any $x' \in \mathbb{R}^n$ there exists y' with $Ay' = x'$; by definition of f_A as the infimum, the preceding display is thus equivalent to

$$f(y') \geq f_A(Ay') \geq f(y) + \langle s, Ay' - Ay \rangle \quad \text{for all } y' \in \mathbb{R}^m.$$

This holds if and only if $A^T s$ is a subgradient of f at y . □

Appendix C

Optimality, stability, and duality

The existence and continuity properties of minimizers of (convex) optimization problems play a central role in much of statistical theory. They are especially essential in our understanding of loss functions and the associated optimality properties. In our context, this is especially central for problems of classification calibration or surrogate risk consistency, as in Chapters 16. This appendix records several representative results along these lines, and also builds up the duality theory associated with convex conjugates, frequently identified as Fenchel-Young duality.

Broadly, throughout this appendix, we shall consider the generic optimization problem

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && x \in C \end{aligned} \tag{C.0.1}$$

where C is a closed convex set (we have not yet assumed convexity of f), Throughout (as in the previous appendix) we assume that f is proper, so that $f(x) > -\infty$ for each x , and that $f(x) = +\infty$ if $x \notin \text{dom } f$.

The most basic question we might ask is when minimizers even exist in the problem (C.0.1). The standard result in this vein is that if minimizers exist whenever C is compact and f is lower semicontinuous (B.3.9), that is, its epigraph is closed, i.e., $\liminf_{x \rightarrow x_0} f(x) \geq f(x_0)$.

Proposition C.0.1. *Let C be compact and $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be lower semi-continuous (B.3.9) over C . Then $\inf_{x \in C} f(x) > -\infty$ and the infimum is attained.*

Proof Let $f^* = \inf_{x \in C} f(x)$, where for now we allow the possibility that $f^* = -\infty$. Let $x_n \in C$ be a sequence of points satisfying $f(x_n) \rightarrow f^*$. Proceeding to a subsequence if necessary, we can assume that $x_n \rightarrow x^* \in C$ by the compactness of C . Then lower semi-continuity guarantees that $f^* = \lim_n f(x_n) \geq f(x^*) \geq f^*$, and so $f(x^*) = f^*$ and so necessarily $f^* > -\infty$. \square

When the domain C is not compact but only closed, alternative conditions are necessary to guarantee the existence of minimizers. Perhaps the most frequent, and one especially useful with convexity (as we shall see), is that f is *coercive*, meaning that

$$f(x) \rightarrow \infty \text{ whenever } \|x\| \rightarrow \infty.$$

Proposition C.0.2. *Let C be closed and $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be lower semi-continuous over C and coercive. Then $\inf_{x \in C} f(x) > -\infty$ and the infimum is attained.*

Proof Once again, let $f^* = \inf_{x \in C} f(x)$ and let $x_n \in C$ satisfy $f(x_n) \rightarrow f^*$. Certainly x_n must be a bounded sequence because f is coercive. Thus, it has a subsequent limit, and w.l.o.g. we assume that $x_n \rightarrow x^* \in C$ by closedness. Lower semi-continuity guarantees that $f^* \geq \liminf_n f(x_n) = f(x^*) \geq f^*$, giving the result. \square

Finally, we make a small remark norms and dual norms, as these will be important for the more quantitative smoothness guarantees we provide. For a norm $\|\cdot\|$, the dual norm $\|\cdot\|_*$ has definition

$$\|y\|_* := \sup_x \{\langle x, y \rangle \mid \|x\| \leq 1\}.$$

This is a norm as it is positively homogeneous, $\|y\|_* = 0$ if and only if $y = 0$, and satisfies the triangle inequality. A few brief examples follow, which we leave as exercises to the reader.

- (i) The ℓ_2 -norm $\|x\|_2 = \sqrt{\langle x, x \rangle}$ is self-dual, so that its dual is $\|\cdot\|_2$.
- (ii) The ℓ_1 and ℓ_∞ norms are dual, that is, $\|x\|_\infty = \sup_{\|y\|_1 \leq 1} \langle x, y \rangle$ and $\|y\|_1 = \sup_{\|x\|_\infty \leq 1} \langle x, y \rangle$.
- (iii) For all $p \in [1, \infty]$, the dual to the ℓ_p norm $\|x\|_p = (\sum_{j=1}^d |x_j|^p)^{1/p}$ is the ℓ_q norm with $q = \frac{p}{p-1}$, that is, for the $q \geq 1$ satisfying $\frac{1}{p} + \frac{1}{q} = 1$.

C.1 Optimality conditions and stability properties

With the basic results on existence of minimizers in place, we turn to convex optimization problems, where f is closed convex and C is a closed convex set, and we assume essentially without loss of generality that $\text{dom } f \supset \text{int } C$ (as otherwise, we may replace C with $C \cap \text{cl dom } f$). The benefits of convexity appear immediately: f has no local but non-global minimizers, and moreover, if f is strictly convex, then any minimizers (if they exist) are unique.

Proposition C.1.1. *Let f be convex. Then if x is a local minimizer of f over C , it is a global minimizer of f over C . If f is strictly convex, then x is unique.*

Proof To say that x is a local minimizer of f over C is to say that $f(x) \leq f(x')$ for all $x' \in C$ with $\|x' - x\| \leq \epsilon$ for some $\epsilon > 0$. Now, consider $y \in C$. By taking $\lambda > 0$ small enough, we have both $(1 - \lambda)x + \lambda y \in C$ and $\|(1 - \lambda)x + \lambda y - x\| \leq \epsilon$, and so

$$f(x) \leq f((1 - \lambda)x + \lambda y) \leq (1 - \lambda)f(x) + \lambda f(y),$$

which rearranged yields $f(y) \geq f(x)$. If f is additionally strictly convex (recall Corollary B.3.2), then the preceding inequality is strict whenever $y \neq x$. \square

C.1.1 Subgradient characterizations for optimality

First-order stationary conditions are sufficient for global optimality in convex problems. We can say more once we consider subgradients:

Observation C.1.2. *Let f be convex and subdifferentiable at x . Then x minimizes f if and only if $0 \in \partial f(x)$.*

Proof If $0 \in \partial f(x)$, then $f(y) \geq f(x) + \langle 0, y - x \rangle = f(x)$ for all y . Conversely, if x minimizes f , then we have $f(y) \geq f(x)$ for all y , and in particular, $0 \in \partial f(x)$. \square

Things become a bit more complicated when we consider the constraints in the problem (C.0.1), so that the point x may be restricted. In this case, it is important and useful to consider the *normal cone* to the set C , which is (essentially) the collection of vectors pointing out of C .

Definition C.1. Let C be a closed convex set. The normal cone to C at the point $x \in C$ is the collection of vectors

$$\mathcal{N}_C(x) := \{v \mid \langle v, y - x \rangle \leq 0 \text{ for all } y \in C\}.$$

So $\mathcal{N}_C(x)$ is the collection of vectors making an obtuse angle with any direction into the set C from x . **JCD Comment:** Draw a picture, and also, put this earlier in the discussion of convex sets.

It is clear that $\mathcal{N}_C(x)$ is indeed a cone: if $v \in \mathcal{N}_C(x)$, then certainly $tv \in \mathcal{N}_C(x)$ for all $t \geq 0$. It is closed convex, being the intersection of halfspaces. Moreover, if $x \in \text{int } C$, then we have $\mathcal{N}_C(x) = \{0\}$, and additionally, we can connect the supporting hyperplanes of C to its normal cones: Theorem B.1.14 gives the following corollary.

Corollary C.1.3. Let C be closed convex. Then for any $x \in \text{bd}(C)$, the normal cone $\mathcal{N}_C(x)$ is non-trivial and consists of the collection of supporting hyperplanes to C at x .

By a bit of subgradient calculus, we can then write optimality conditions involving the normal cones to C . If C is a closed convex set, the convex indicator function $\mathbf{I}_C(x)$ has subdifferentials

$$\partial \mathbf{I}_C(x) = \begin{cases} \{0\} & \text{if } x \in \text{int } C \\ \mathcal{N}_C(x) & \text{if } x \in \text{bd}(C) \\ \emptyset & \text{otherwise.} \end{cases}$$

The only case requiring justification is the boundary case; for this, we note that $w \in \mathcal{N}_C(x)$ if and only if $\langle w, y - x \rangle \leq 0$ for all $y \in C$, which in turn occurs if and only if $\mathbf{I}_C(y) \geq \mathbf{I}_C(x) + \langle w, y - x \rangle$ for all y .

The subdifferential calculation for $\mathbf{I}_C(x)$ yields the following general optimality characterization for problem (C.0.1).

Proposition C.1.4. In the problem (C.0.1), let $x \in \text{int dom } f$. Then x minimizes f over C if and only if

$$0 \in \partial f(x) + \mathcal{N}_C(x). \quad (\text{C.1.1})$$

Proof The minimization problem (C.0.1) is equivalent to the problem

$$\underset{x}{\text{minimize}} \quad f(x) + \mathbf{I}_C(x).$$

As $x \in \text{int dom } f$, f has nonempty compact convex subdifferential $\partial f(x)$, and so $\partial(f + \mathbf{I}_C)(x) = \partial f(x) + \partial \mathbf{I}_C(x) = \partial f(x) + \mathcal{N}_C(x)$ by Proposition B.3.25. Apply Observation C.1.2. \square

Several equivalent versions of Proposition C.1.4 are possible. The first is that

$$-\partial f(x) \cap \mathcal{N}_C(x) \neq \emptyset,$$

that is, there is a subgradient vector $g \in \partial f(x)$ such that $-g \in \mathcal{N}_C(x)$, so that $-g$ points outside the set C . **JCD Comment:** Draw a picture

Another variant, frequently used, is to write Proposition C.1.4 as that x solves problem (C.0.1) if and only if there exists $g \in \partial f(x)$ such that

$$\langle g, y - x \rangle \geq 0 \quad \text{for all } y \in C. \quad (\text{C.1.2})$$

Indeed, taking $g \in \partial f(x)$ to be the element satisfying $-g \in \mathcal{N}_C(x)$, we immediately see that $\langle -g, y - x \rangle \leq 0$ for all $y \in C$ by definition of the normal cone, which is (C.1.2). **JCD Comment:** Use picture above

C.1.2 Stability properties of minimizers

The characterizations (C.1.1) and (C.1.2) of optimality for convex optimization problems allow us to develop some stability properties of the minimizers of convex problems. These, in turn, will relate to smoothness properties of various dual functions (as we explore in the sequel), which again become important in the study of consistent losses. Here, we collect a few of the typical results. Typical results in this vein exhibit a few properties: that solutions are “stable,” meaning that small tilts of the function f do not change solutions significantly, or that the function f exhibits various strong growth properties.

JCD Comment: Add a bit of commentary about when there is growth then we expect minimizers to be smoothish.

For our starting point, we begin by consider *strongly convex* functions, where a function is λ -strongly convex with respect to the norm $\|\cdot\|$ if for all $t \in [0, 1]$ and $x, y \in \text{dom } f$,

$$f(tx + (1-t)y) \leq tf(x) + (1-t)f(y) - \frac{\lambda}{2}t(1-t)\|x - y\|^2. \quad (\text{C.1.3})$$

The definition (C.1.3) makes strict convexity quantitative in a fairly precise way, and has several equivalent characterizations.

Proposition C.1.5 (Equivalent characterizations of strong convexity). *Let f be a convex function, subdifferentiable on its domain. Then the following are equivalent.*

(i) f is λ -strongly convex (Eq. (C.1.3)).

(ii) For all $y \in \mathbb{R}^d$ and $g \in \partial f(x)$,

$$f(y) \geq f(x) + \langle g, y - x \rangle + \frac{\lambda}{2}\|x - y\|^2.$$

(iii) For all $x, y \in \text{dom } f$ and $g_x \in \partial f(x)$ and $g_y \in \partial f(y)$,

$$\langle g_x - g_y, x - y \rangle \geq \lambda\|x - y\|^2.$$

Proof Let us prove that inequality (ii) holds if and only if (iii) does. Let $g_x \in \partial f(x)$ and $g_y \in \partial f(y)$ and assume (ii) holds. Then

$$\begin{aligned} f(y) &\geq f(x) + \langle g_x, y - x \rangle + \frac{\lambda}{2}\|y - x\|^2 \\ f(x) &\geq f(y) + \langle g_y, x - y \rangle + \frac{\lambda}{2}\|x - y\|^2 \end{aligned}$$

and adding the equations we obtain

$$0 \geq \langle g_x - g_y, y - x \rangle + \lambda \|x - y\|^2.$$

Rearranging gives part (iii). Conversely, assume (iii), and for $t \in [0, 1]$ let $x_t = (1 - t)x + ty$ and define $h(t) = f(x_t)$. Then h is convex and hence almost everywhere differentiable (and locally Lipschitz), so that $h(1) = h(0) + \int_0^1 h'(t)dt$. Noting that

$$h'(t) = \langle g_t, y - x \rangle \text{ for some } g_t \in \partial f(x_t)$$

(recall the subgradient characterization of Proposition B.3.20), we have

$$h'(t) = \langle g_t, y - x \rangle = \langle g_t - g_x, y - x \rangle + \langle g_x, y - x \rangle = \frac{1}{t} \langle g_t - g_x, (1 - t)x + ty - x \rangle + \langle g_x, y - x \rangle$$

and so as $h(1) = f(y)$ and $h(0) = f(x)$,

$$\begin{aligned} f(y) &= h(0) + \int_0^1 \frac{\langle g_t - g_x, (1 - t)x + ty - x \rangle}{t} dt + \langle g_x, y - x \rangle \\ &\geq f(x) + \int_0^1 \frac{\lambda \|(1 - t)x + ty - x\|^2}{t} dt + \langle g_x, y - x \rangle \\ &= f(x) + \lambda \|x - y\|^2 \int_0^1 t dt + \langle g_x, y - x \rangle = f(x) + \langle g_x, y - x \rangle + \frac{\lambda}{2} \|y - x\|^2. \end{aligned}$$

That (ii) implies (i) is relatively straightforward: we have

$$\begin{aligned} f(y) &\geq f(tx + (1 - t)y) + t \langle g_t, y - x \rangle + \frac{\lambda}{2} t^2 \|x - y\|^2 \\ f(x) &\geq f(tx + (1 - t)y) + (1 - t) \langle g_t, x - y \rangle + \frac{\lambda}{2} (1 - t)^2 \|x - y\|^2 \end{aligned}$$

for any $g_t \in \partial f(tx + (1 - t)y)$. Multiply the first inequality by $(1 - t)$ and the second by t , then add them to obtain

$$tf(x) + (1 - t)f(y) \geq f(tx + (1 - t)y) + \frac{\lambda}{2} [(1 - t)t^2 + t(1 - t)^2] \|x - y\|^2,$$

and note that $(1 - t)t^2 + t(1 - t)^2 = t(1 - t)$. Finally, let (i) hold, and which is equivalent to the condition that

$$\frac{f((1 - t)x + ty) - f(x)}{t} + \frac{\lambda}{2} (1 - t) \|x - y\|^2 \leq f(y) - f(x)$$

for $t \in (0, 1)$. Taking $t \downarrow 0$ gives $f'(x; y - x) + \frac{\lambda}{2} \|x - y\|^2 \leq f(y) - f(x)$, and because $f'(x; y - x) = \sup_{s \in \partial f(x)} \langle s, y - x \rangle$ we obtain (ii). \square

As a first example application of strong convexity, consider minimizers of the tilted functions

$$f_u(x) := f(x) - \langle u, x \rangle$$

as u varies. First, note that minimizers necessarily exist: the function $f_u(x) \rightarrow \infty$ whenever $\|x\| \rightarrow \infty$ by condition (ii) in Proposition C.1.5, and so we can restrict to minimizing f_u over

compacta. Moreover, the minimizers $x_u := \operatorname{argmin}_x f_u(x)$ are unique, as the functions f_u are strongly (and hence strictly) convex. However, we can say more. Indeed, let C be any closed convex set and let

$$x_u = \operatorname{argmin}_{x \in C} f_u(x). \quad (\text{C.1.4})$$

We claim the following:

Proposition C.1.6. *Let f be λ -strongly convex with respect to the norm $\|\cdot\|$ and subdifferentiable on C . Then the mapping $u \mapsto x_u$ is $\frac{1}{\lambda}$ -Lipschitz continuous with respect to the dual norm $\|\cdot\|_*$, that is, $\|x_u - x_v\| \leq \frac{1}{\lambda} \|u - v\|_*$.*

Proof We use the optimality condition (C.1.2). We have $\partial f_u(x) = \partial f(x) - u$, and thus for any u, v we have both

$$\langle g_u - u, y - x_u \rangle \geq 0 \quad \text{and} \quad \langle g_v - v, y - x_v \rangle \geq 0$$

for some $g_u \in \partial f(x_u)$ and $g_v \in \partial f(x_v)$ for all $y \in C$. Set $y = x_v$ in the first inequality and $y = x_u$ in the second and add them to obtain

$$\langle g_u - g_v + v - u, x_v - x_u \rangle \geq 0 \quad \text{or} \quad \langle v - u, x_v - x_u \rangle \geq \langle g_v - g_u, x_v - x_u \rangle.$$

By strong convexity the last term satisfies $\langle g_v - g_u, x_v - x_u \rangle \geq \lambda \|x_u - x_v\|^2$. By definition of the dual norm, $\|v - u\|_* \|x_v - x_u\| \geq \langle v - u, x_v - x_u \rangle$, so $\|u - v\|_* \|x_v - x_u\| \geq \lambda \|x_u - x_v\|^2$, which is the desired result. \square

JCD Comment: Figure for the preceding lemma.

There are alternative versions of strong convexity, typically given the name *uniform convexity* in the convex analysis literature, which allow generalizations and similar quantitative stability properties. In analogy with the strong convexity condition (C.1.3), we say that f is (λ, κ) -uniformly convex, where $\kappa \geq 2$, over C if it is closed and for all $t \in [0, 1]$ and $x, y \in C$,

$$f(tx + (1-t)y) \leq tf(x) + (1-t)f(y) - \frac{\lambda}{2} t(1-t) \|x - y\|^\kappa [(1-t)^{\kappa-1} + t^{\kappa-1}]. \quad (\text{C.1.5})$$

Notably, the $\kappa = 2$ case is simply the familiar strong convexity property. Taking $\kappa > 2$ weakens strong convexity, yielding correspondingly weaker guarantees of stability and optimality. We can, however, provide analogies to Propositions C.1.5 and C.1.6.

Proposition C.1.7 (Equivalent characterizations of uniform convexity). *Let f be a convex function, subdifferentiable on its domain. Then the following are equivalent:*

- (i) f is (λ, κ) -uniformly convex with respect to the norm $\|\cdot\|$.
- (ii) For all $y \in \mathbb{R}^d$ and $g \in \partial f(x)$,

$$f(y) \geq f(x) + \langle g, y - x \rangle + \frac{\lambda}{2} \|x - y\|^\kappa.$$

Additionally, either of (i) or (ii) imply that

(iii) For all $x, y \in \text{dom } f$ and $g_x \in \partial f(x)$ and $g_y \in \partial f(y)$,

$$\langle g_x - g_y, x - y \rangle \geq \lambda \|x - y\|^\kappa,$$

which in turn implies that f is $(\frac{2}{\kappa}\lambda, \kappa)$ -uniformly convex with respect to the norm $\|\cdot\|$.

We leave the proof of the proposition as Exercise C.2, noting that it follows via the same arguments as those we use to prove the strong convexity version (Proposition C.1.5). An analog of Proposition C.1.6 also holds with an alternative smoothness condition.

Proposition C.1.8. *Let f be (λ, κ) -uniformly convex with respect to the norm $\|\cdot\|$ and subdifferentiable on C . Then the mapping $u \mapsto x_u$ is $\frac{1}{\kappa-1}$ -Hölder, and in particular,*

$$\|x_u - x_v\| \leq \lambda^{-\frac{1}{\kappa-1}} \|u - v\|_*^{\frac{1}{\kappa-1}}$$

for all u, v .

Proposition C.1.8 follows as an immediate consequence of Proposition C.2.7 to come, where we connect smoothness of the minimizers x_u with differentiability properties of the conjugate function.

JCD Comment: Add some figures on strict convexity implying a bit of growth around a neighborhood, and stability properties of strongly convex functions.

The weakest version of such strong convexity properties is strict convexity, for which a careful reading of the proof of Proposition C.1.5 (replace all λ with 0 and inequalities with strict inequalities) gives the following characterization of equivalent definitions of strict convexity (recall also Corollary B.3.2).

Corollary C.1.9. *Let f be a convex function subdifferentiable on C . The following are equivalent.*

- (i) f is strictly convex on C .
- (ii) For all $x \in C$, $y \neq x$, and $g \in \partial f(x)$,

$$f(y) > f(x) + \langle g, y - x \rangle.$$

- (iii) For all $x, y \in C$ and $g_x \in \partial f(x)$ and $g_y \in \partial f(y)$,

$$\langle g_x - g_y, x - y \rangle > 0.$$

Using Corollary C.1.9, we can then obtain certain smoothness properties of the tilted minimizers x_u of the minimization (C.1.4). We begin with a lemma that guarantees growth of convex functions over their first-order approximations.

Lemma C.1.10. *Let f be convex and subdifferentiable on the closed convex set C , and for any fixed $g \in \partial f(x_0)$ define the Bregman divergence*

$$D(x, x_0) := f(x) - f(x_0) - \langle g, x - x_0 \rangle.$$

Then for all $0 \leq \epsilon \leq \epsilon'$, $\delta(\epsilon) := \inf\{D(x, x_0) \mid x \in C, \|x - x_0\| \geq \epsilon\}$ is attained, nonnegative, and $\delta(\epsilon') \geq \frac{\epsilon'}{\epsilon} \delta(\epsilon)$.

Proof Fix $x \in C$. Letting $h(t) = D(x_0 + t(x - x_0), x_0)$, h is convex in $t \geq 0$, locally Lipschitz, and satisfies $h(0) = \inf_t h(t) = 0$, so we can write $h(t) = h(0) + \int_0^t h'(s; 1) ds$. Additionally, $s \mapsto h'(s; 1) \geq 0$ is nondecreasing by the increasing slopes criterion (B.3.4).

For all $\epsilon > 0$, then, we may restrict infimum in the definition of $\delta(\epsilon)$ to those $x \in C$ satisfying $\|x - x_0\| = \epsilon$, a compact set, so that the infimum is attained at some $x_\epsilon \in C$ with $\|x_\epsilon - x_0\| = \epsilon$. Now, let $\epsilon' > \epsilon$, and $x_{\epsilon'}$ achieve the infimum in $\delta(\epsilon)$. Then setting $x' = \frac{\epsilon'}{\epsilon}(x_{\epsilon'} - x_0) + x_0$ (implying $x_{\epsilon'} = \frac{\epsilon'}{\epsilon}(x' - x_0) + x_0$), we have $\|x' - x_0\| = \epsilon$ and so $D(x', x_0) \geq D(x_\epsilon, x_0) = \delta(\epsilon)$. Set $h(t) = D(x_0 + t(x' - x_0), x_0)$. Rewriting and using the first-order convexity condition,

$$\begin{aligned} \delta(\epsilon') = D(x_{\epsilon'}, x_0) &= D\left(\frac{\epsilon'}{\epsilon}(x' - x_0) + x_0, x_0\right) = h\left(\frac{\epsilon'}{\epsilon}\right) \geq h(1) + \left[\frac{\epsilon'}{\epsilon} - 1\right] h'(1; 1) \\ &= D(x', x_0) + \left[\frac{\epsilon'}{\epsilon} - 1\right] h'(1; 1). \end{aligned}$$

A minor variant of the criterion of increasing slopes (B.3.4) and that $h(0) = 0$ then gives $h'(1; 1) = \lim_{t \downarrow 0} \frac{h(1+t) - h(1)}{t} \geq \frac{h(1) - h(0)}{1} = h(1) = D(x', x_0)$, so we have

$$\delta(\epsilon') = D(x_{\epsilon'}, x_0) \geq D(x', x_0) + \left[\frac{\epsilon'}{\epsilon} - 1\right] D(x', x_0) = \frac{\epsilon'}{\epsilon} D(x', x_0) \geq \frac{\epsilon'}{\epsilon} \delta(\epsilon)$$

as desired. \square

Whenever f is strictly convex, because the infimum in $\delta(\epsilon)$ is attained in Lemma C.1.10, we have the following guarantee.

Lemma C.1.11. *Let the conditions of Lemma C.1.10 hold and additionally let f be strictly convex. Then $\delta(\epsilon) > 0$ for all $\epsilon > 0$.*

Combining these results yields the following non-quantitative version of Proposition C.1.6:

Proposition C.1.12. *Let f be strictly convex and subdifferentiable on the closed convex set C , and assume that the minimum $x_0 = \operatorname{argmin}_{x \in C} f(x)$ is attained. Then the mapping $u \mapsto x_u$ is continuous in a neighborhood of $u = 0$.*

Proof We show first that x_u is continuous at $u = 0$. By Lemmas C.1.10 and C.1.11, we see that for $x \in C$ we have

$$f(x) \geq f(x_0) + \langle g, x - x_0 \rangle + \delta(\|x - x_0\|) \geq f(x_0) + \delta(\|x - x_0\|),$$

where $g \in \partial f(x_0)$ satisfies $\langle g, x - x_0 \rangle \geq 0$ for all $x \in C$ by the optimality condition (C.1.2). Now, pick $\epsilon > 0$, so that if $\|x - x_0\| > \epsilon$ we have $\delta(\|x - x_0\|) \geq \|x - x_0\| \frac{\delta(\epsilon)}{\epsilon}$ by Lemma C.1.10. Then if u satisfies $\|u\| < \frac{\delta(\epsilon)}{\epsilon}$, we have

$$\begin{aligned} f(x) - \langle u, x \rangle &= f(x) - \langle u, x - x_0 \rangle - \langle u, x_0 \rangle \\ &\geq f(x_0) - \langle u, x_0 \rangle + \delta(\|x - x_0\|) - \langle u, x - x_0 \rangle \\ &\geq f(x_0) - \langle u, x_0 \rangle + \delta(\|x - x_0\|) - \langle u, x - x_0 \rangle \\ &\geq f(x_0) - \langle u, x_0 \rangle + \delta(\|x - x_0\|) - \|u\| \|x - x_0\| \\ &> f(x_0) - \langle u, x_0 \rangle + \frac{\delta(\epsilon)}{\epsilon} \|x - x_0\| - \frac{\delta(\epsilon)}{\epsilon} \|x - x_0\| = f(x_0) - \langle u, x_0 \rangle. \end{aligned}$$

Thus any minimizer x_u of $f(x) - \langle u, x \rangle$ over $x \in C$ must satisfy $\|x_u - x_0\| \leq \epsilon$, and strict convexity guarantees its uniqueness.

The argument that $u \mapsto x_u$ is continuous in a neighborhood of zero is completely similar once we recognize that for the divergence $D_f(x, x_0) := f(x) - \langle g, x - x_0 \rangle - f(x_0)$ (where $g \in \partial f(x_0)$ is fixed), we have $D_f = D_{f_u}$ for $f_u(x) = f(x) - \langle u, x \rangle$ and x_u is near x_0 for u small. \square

C.2 Conjugacy and duality properties

Attached to any function is its *convex conjugate*, sometimes called the *Fenchel* or *Fenchel-Legendre* conjugate function, defined by

$$f^*(s) := \sup_x \{ \langle s, x \rangle - f(x) \}. \quad (\text{C.2.1})$$

For any f , the conjugate f^* is a closed convex function, as it is the supremum of linear functions. This function helps to exhibit a duality for convex functions similar to those for convex sets, which we can describe as the intersection of all halfspaces containing them (recall Theorem B.1.15 and the equalities (B.1.3)–(B.1.4)).

JCD Comment: Draw a picture of the conjugate

The conjugate function is the largest gap between the linear functional $x \mapsto \langle s, x \rangle$ and the function f itself. The remarkable property of such conjugates is that their biconjugates describe the function f itself, or at least the largest closed convex function below f . To make this a bit more precise, we state a theorem, and then connect to so-called *convex closures* of functions.

Theorem C.2.1. *Let f be closed convex and f^* be its conjugate (C.2.1). Then*

$$f^{**}(x) = f(x) \text{ for all } x.$$

Proof By definition, we have

$$f^{**}(x) = \sup_s \{ \langle x, s \rangle - f^*(s) \},$$

and we always have $\langle x, s \rangle - f^*(s) \leq f(x)$ by definition of $f^*(s) = \sup_x \{ \langle s, x \rangle - f(x) \}$. So immediately we see that $f^{**}(x) \leq f(x)$.

We essentially show that the linear functions $h_s(x) := \langle x, s \rangle - f^*(s)$ describe (enough) of the global linear underestimators of f so that $f(x) = \sup_s h_s(x)$, allowing us to apply Theorem B.3.7. Indeed, let $l(x) = \langle s, x \rangle + b$ be any global underestimator of f . Then we must have $b \leq f(x) - \langle s, x \rangle$ for all x , that is, $b \leq \inf_x \{ f(x) - \langle s, x \rangle \} = -\sup_x \{ \langle s, x \rangle - f(x) \} = -f^*(s)$, that is, $l(x) \leq \langle s, x \rangle - f^*(s) = h_s(x)$. Apply Theorem B.3.7. \square

We may visualize f^{**} as pulling a string up below a function f , yielding the largest closed convex underestimator of f . Combining Theorem C.2.1 with Proposition B.3.14, we obtain the following corollary. (Note that we require $f(x) > -\infty$ for all x .)

Corollary C.2.2. *Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be convex and f^* its conjugate (C.2.1). Then*

$$f^{**}(x) = \text{cl } f(x) \text{ for all } x.$$

*In particular, f is lower semicontinuous at x if and only if $f(x) = \text{cl } f(x) = f^{**}(x)$.*

C.2.1 Gradient dualities and the Fenchel-Young inequality

It is immediate from the definition that for any pair s, x we have the *Fenchel-Young inequality*

$$\langle s, x \rangle \leq f^*(s) + f(x). \quad (\text{C.2.2})$$

Even more, combining Theorem C.2.1 with this observation, we can exhibit a duality between subgradients of f and f^* with this inequality.

Proposition C.2.3. *Let f be closed convex. Then*

$$\langle s, x \rangle = f^*(s) + f(x) \quad \text{if and only if} \quad s \in \partial f(x) \quad \text{if and only if} \quad x \in \partial f^*(s).$$

Proof If $\langle s, x \rangle = f^*(s) + f(x)$, then $-f(x) + \langle s, x \rangle = f^*(s) \geq \langle s, y \rangle - f(y)$ for all y , and rearranging, we have $f(y) \geq f(x) + \langle s, y - x \rangle$, that is, $s \in \partial f(x)$. Conversely, if $s \in \partial f(x)$ then $0 \in \partial f(x) - s$, so that x minimizes $f(x) - \langle s, x \rangle$, or equivalently, x maximizes $\langle s, x \rangle - f(x)$ and so $\langle s, x \rangle - f(x) = \sup_x \{\langle s, x \rangle - f(x)\}$ as desired. The final statement is immediate from a parallel argument and the duality in Theorem C.2.1. \square

Writing Proposition C.2.3 differently, we see that ∂f and ∂f^* are inverses of one another. That is, as set-valued mappings, where

$$(\partial f)^{-1}(s) := \{x \mid s \in \partial f(x)\},$$

we have the following corollary.

Corollary C.2.4. *Let f and f^* be subdifferentiable. Then*

$$\partial f^* = (\partial f)^{-1} \quad \text{and} \quad \partial f = (\partial f^*)^{-1}$$

and

$$\partial f^*(s) = \operatorname{argmax}_x \{\langle s, x \rangle - f(x)\} \quad \text{and} \quad \partial f(x) = \operatorname{argmax}_s \{\langle s, x \rangle - f^*(s)\}.$$

Notably, if f and f^* are differentiable, then $\nabla f = (\nabla f^*)^{-1}$.

Additionally, we see that the domains and images of ∂f and ∂f^* are also related, which guarantees convexity properties of their images as well.

Corollary C.2.5. *Let f be closed convex. Then*

$$\operatorname{dom} \partial f = \operatorname{Im} \partial f^* \quad \text{and} \quad \operatorname{dom} \partial f^* = \operatorname{Im} \partial f.$$

Proof Let $x \in \operatorname{dom} \partial f$, so that $\partial f(x)$ is non-empty. Then $s \in \partial f(x)$ implies that $\langle s, x \rangle = f(x) + f^*(s)$ and $x \in \partial f^*(s)$ by Proposition C.2.3. Similarly, if $x \in \operatorname{Im} \partial f^*$, then there is some s for which $x \in \partial f^*(s)$ and so $\langle s, x \rangle = f(x) + f^*(s)$ and $s \in \partial f(x)$. \square

C.2.2 Smoothness and strict convexity of conjugates

The dualities in derivative mappings extend to various smoothness dualities, which can be quite useful as well. These types of results build from the stability properties of solution mappings, as in those for tilted minimizers (C.1.4) in Propositions C.1.6 and C.1.12. They also relate different smoothness properties of f and f^* , as well as their domains of definition, to the existence and continuity of minimizers of $f(x) - \langle s, x \rangle$.

When we assume that f has quantitative strong convexity or smoothness properties, we can give similar quantitative guarantees for the smoothness and strong convexity of f^* :

Proposition C.2.6. *Let $f : \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ be λ -strongly convex with respect to the norm $\|\cdot\|$ (see Eq. (C.1.3)) on its domain. Then $\text{dom } f^* = \mathbb{R}^d$ and ∇f^* is $\frac{1}{\lambda}$ -Lipschitz continuous with respect to the dual norm $\|\cdot\|_*$, that is,*

$$\|\nabla f^*(u) - \nabla f^*(v)\| \leq \frac{1}{\lambda} \|u - v\|_*$$

for all u, v . Conversely, let $f : \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ be convex with L -Lipschitz gradient with respect to $\|\cdot\|$ on \mathbb{R}^d . Then f^* is $\frac{1}{L}$ -strongly convex with respect to the dual norm $\|\cdot\|_*$ on convex subsets $C \subset \text{dom } \partial f^*$, and in particular,

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|_*^2. \quad (\text{C.2.3})$$

Proof For the first claim, let $C = \text{dom } f$. Then Proposition C.1.6 shows that if $x_1 = \text{argmin}_x \{f(x) - \langle s_1, x \rangle\}$ and $x_2 = \text{argmin}_x \{f(x) - \langle s_2, x \rangle\}$ (which exist and are necessarily unique), we have $\|x_1 - x_2\| \leq \frac{1}{\lambda} \|s_1 - s_2\|_*$. Then Proposition C.2.3 shows that $x_i \in \partial f^*(s_i)$ for $i = 1, 2$, and hence $\partial f^*(s_i)$ is necessarily single-valued and $(1/\lambda)$ -Lipschitz continuous.

The converse is a bit trickier. Let x and y be arbitrary and $s_x = \nabla f(x)$ and $s_y = \nabla f(y)$; we prove inequality (C.2.3), known as *co-coercivity*. By the L -Lipschitz continuity of ∇f , we have

$$\begin{aligned} f(y) &= f(x) + \int_0^1 \langle \nabla f(x + t(y - x)), y - x \rangle dt \\ &= f(x) + \langle \nabla f(x), y - x \rangle + \int_0^1 \langle \nabla f(x + t(y - x)) - \nabla f(x), y - x \rangle dt \\ &\leq f(x) + \langle s_x, y - x \rangle + \int_0^1 L t \|y - x\|^2 dt = f(x) + \langle s_x, y - x \rangle + \frac{L}{2} \|y - x\|^2, \end{aligned}$$

which is valid for any x, y . Note that $f(x) - \langle s_x, x \rangle = -f^*(s_x)$, so that rearranging we have

$$\begin{aligned} f^*(s_x) &\leq \langle s_x, y \rangle - f(y) + \frac{L}{2} \|y - x\|^2 = \langle s, y \rangle - f(y) + \langle s_x - s, y \rangle + \frac{L}{2} \|y - x\|^2 \\ &\leq f^*(s) + \langle s_x - s, y \rangle + \frac{L}{2} \|y - x\|^2, \end{aligned}$$

valid for any vector s and any y . We may in particular take an infimum over y on the right hand side, where

$$\begin{aligned} \inf_y \langle s_x - s, y \rangle + \frac{L}{2} \|y - x\|^2 &= \inf_y \langle s_x - s, y - x \rangle + \frac{L}{2} \|y - x\|^2 + \langle s_x - s, x \rangle \\ &\stackrel{(*)}{=} \inf_t \left\{ t \|s_x - s\|_* + \frac{Lt^2}{2} \right\} + \langle s_x - s, x \rangle \\ &= -\frac{1}{2L} \|s_x - s\|_*^2 + \langle s_x - s, x \rangle, \end{aligned}$$

where equality (\star) follows by definition of the dual norm and we identify $t = \|y - x\|$. Thus

$$f^*(s_x) + \langle x, s - s_x \rangle + \frac{1}{2L} \|s - s_x\|_*^2 \leq f^*(s)$$

for all s . As $x \in \partial f^*(s_x)$, Proposition C.1.5(ii) gives the strong convexity result. The rest is algebraic manipulations with $s_y = \nabla f(y)$ and an application of Proposition C.1.5, part (iii). \square

We can extend Proposition C.2.6 to the uniformly convex case (recall Proposition C.1.7, though for this we require a slightly extended definition of smoothness beyond Lipschitz continuity of the gradients.

Definition C.2. A function f is $(L, \beta, \|\cdot\|)$ -smooth if it has β -Hölder continuous gradient with respect to the norm $\|\cdot\|$, that is,

$$\|\nabla f(x) - \nabla f(y)\|_* \leq L \|x - y\|^\beta \quad \text{for } x, y \in \text{dom } f.$$

We then have the following analog of Proposition C.2.6.

Proposition C.2.7. Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be (λ, κ) -uniformly convex (C.1.5) with respect to the norm $\|\cdot\|$. Then $\text{dom } f^* = \mathbb{R}^d$ and ∇f^* is $(\lambda^{-\frac{1}{\kappa-1}}, \frac{1}{\kappa-1}, \|\cdot\|_*)$ -smooth. Conversely, let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and $(L, \beta, \|\cdot\|)$ -smooth. Then f^* is $(2^{\frac{\beta-1}{\beta}} L^{-\frac{1}{\beta-1}}, \frac{\beta}{\beta-1})$ -uniformly convex (C.1.5) with respect to the dual norm $\|\cdot\|_*$ on any convex subset of $\text{dom } \partial f^*$, and in particular,

$$f^*(s_1) \geq f^*(s_0) + \langle \partial f^*(s_0), s_1 - s_0 \rangle + \frac{\beta-1}{\beta} L^{-\frac{1}{\beta-1}} \|s_0 - s_1\|_*^{\frac{\beta}{\beta-1}}.$$

Exercise C.3 asks for a proof of this proposition, which more or less follows from a similar technique as that we use to prove Proposition C.2.6. In short, however, uniform convexity and smoothness are dual to one another (with a minor loss in leading constant multipliers): let κ and $\kappa^* = \frac{\kappa}{\kappa-1}$ be conjugates, where $\kappa \geq 2$. Then the function f is uniformly convex with exponent κ if and only if f^* has κ^* -Hölder continuous gradients.

There are more qualitative versions of Proposition C.2.6 that allow us to give a duality between strict convexity and continuous differentiability of f . Here, we give one typical result.

Proposition C.2.8. Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be strictly convex and closed. Then $\text{int dom } f^* \neq \emptyset$ and f^* is continuously differentiable on $\text{int dom } f^*$. Conversely, let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be differentiable on $\Omega := \text{int dom } f$. Then f^* is strictly convex on each convex $C \subset \nabla f(\Omega)$.

These results should be roughly expected because of the duality that $\nabla f = (\nabla f^*)^{-1}$ and that $\partial f^*(s) = \text{argmin}_x \{\langle s, x \rangle - f(x)\}$, because strict convexity guarantees uniqueness of minimizers (Proposition C.1.1) so that ∂f^* should be a singleton.

Proof To see that $\text{int dom } f^*$ is non-empty, we use the identification $f'_\infty(v) = \sigma_{\text{dom } f^*}(v)$ in Proposition C.3.5 and the interior identification in Proposition B.2.7. Because f is strictly convex, for any $x \in \text{dom } f$ we have

$$0 < \frac{f(x - tv) - f(x)}{t} + \frac{f(x + td) - f(x)}{t} \quad \text{for } t > 0,$$

and taking $t \rightarrow \infty$ gives $0 < f'_\infty(-v) + f'_\infty(v)$. Proposition B.2.7 then shows that $\text{int dom } f^* \neq \emptyset$.

For the claim that f^* is continuously differentiable, take $s \in \text{int dom } f^*$, and suppose for the sake of contradiction that $\partial f^*(s)$ has distinct points x_1, x_2 . Then Corollary C.2.4 gives that x_1 and x_2 both minimize $f(x) - \langle s, x \rangle$ over x . But Proposition C.1.1 guarantees $x_1 = x_2$, so that $\partial f^*(s) = \{\nabla f^*(s)\}$ is a singleton, and hence f^* is continuous differentiable at s (Proposition B.3.24).

For the converse claim, let C be a convex set as stated. Suppose for the sake of contradiction that f^* is not strictly convex on C , so that there are distinct points $s_1, s_2 \in C$ for which f^* is affine on the line segment $[s_1, s_2] = \{ts_1 + (1-t)s_2 \mid t \in [0, 1]\}$. As $C \subset \nabla f(\Omega)$ is convex, the midpoint $s = \frac{1}{2}(s_1 + s_2) \in C$ and there exists x satisfying $\nabla f(x) = s$, or $x \in \partial f^*(s)$. Then because f^* is assumed affine on $[s_1, s_2]$, we have $f^*(s) = \frac{1}{2}f^*(s_1) + \frac{1}{2}f^*(s_2)$ and $\langle s, x \rangle = \frac{1}{2}\langle s_1 + s_2, x \rangle$, so

$$\begin{aligned} 0 &= f(x) + f^*(s) - \langle s, x \rangle \\ &= \frac{1}{2}[(f(x) + f^*(s_1) - \langle s_1, x \rangle) + (f(x) + f^*(s_2) - \langle s_2, x \rangle)]. \end{aligned}$$

Each of the terms in parenthesis is 0 if and only if $s_i \in \partial f(x)$, but by assumption $\partial f(x) = \{\nabla f(x)\}$ is a singleton, and we must have $s_1 = s_2$. \square

JCD Comment: Better transition

C.2.3 Smooth convex functions

We close this section by investigating particularly nice classes of functions f , where f and its conjugate f^* are strictly convex and smooth. These results are central to the various conjugate linkage dualities we explore in Chapter 14.3. We therefore make the following definition:

Definition C.3. Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be closed convex. Then f is of Legendre type if

- (i) $\text{int dom } f \neq \emptyset$
- (ii) f is continuously differentiable on $\text{int dom } f$
- (iii) f is strictly convex
- (iv) f satisfies the gradient boundary conditions

$$\|\nabla f(x)\| \rightarrow \infty \text{ as } x \rightarrow \text{bd dom } f \text{ or } \|x\| \rightarrow \infty. \quad (\text{C.2.4})$$

Thus, at the boundaries of their domains or as their argument tends off to infinity, functions of Legendre type have slopes tending to ∞ . This does not guarantee that $f(x) \rightarrow \infty$ as $x \rightarrow \text{bd dom } f$, though it does provide guarantees of regularity that the next theorem highlights.

Theorem C.2.9. Let f be a convex function of Legendre type (Def. C.3). Then f^* is strictly convex, continuously differentiable, and $\text{dom } f^* = \mathbb{R}^d$.

The theorem implies a number of results on continuity of minimizers and tilted minimizers (C.1.4), clarifying some of our earlier results. For example, we have the following corollary.

Corollary C.2.10. Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be a convex function of Legendre type. Then the tilted minimizer

$$x_u := \operatorname{argmin}\{f(x) - \langle u, x \rangle\}$$

exists for all u , is continuous in u and unique, and $x_u \in \text{int dom } f$.

We turn to the proof of the theorem.

Proof of Theorem C.2.9 We state an intermediate lemma, whose proof we defer.

Lemma C.2.11. *Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be closed convex and satisfy the gradient boundary condition that $\|s_n\| \rightarrow \infty$ for any sequence $x_n \rightarrow \text{bd dom } f$ and $s_n \in \partial f(x_n)$. Then*

$$f'_\infty(v) = \infty \text{ for all } v \neq 0$$

if and only if

$$\|s_n\| \rightarrow \infty \text{ whenever } \|x_n\| \rightarrow \infty \text{ and } s_n \in \partial f(x_n).$$

The theorem follows straightforwardly from Lemma C.2.11. By the boundary conditions (C.2.4) associated with f , we have $f'_\infty(v) = \infty$ for all $v \neq 0$ (Lemma C.2.11). Because the support function of $\text{dom } f^*$ satisfies $\sigma_{\text{dom } f^*} = f'_\infty$ (Proposition C.3.5), we see that $\text{dom } f^* = \mathbb{R}^d$ as $\text{dom } f^* = \{s \mid \langle s, v \rangle \leq \sigma_{\text{dom } f^*}(v) \text{ for all } v\}$ (e.g., Proposition B.2.7 or Corollary B.2.2). With this, f^* is continuously differentiable and strictly convex on its domain (Proposition C.2.8). \square

Proof of Lemma C.2.11 As $\text{dom } f^* = \mathbb{R}^d$ if and only if $f'_\infty(v) = \infty$ for all $v \neq 0$ (Corollary C.3.7), it suffices to show the result that $\text{int dom } f^* \neq \mathbb{R}^d$ if and only if there exists an unbounded sequence x_n with $s_n \in \partial f(x_n)$ and for which s_n is convergent.

Let us begin with the unbounded sequence x_n for which $s_n \rightarrow s \in \mathbb{R}^d$; assume for the sake of contradiction that $s \in \text{int dom } f^*$. Because $s_n \in \partial f(x_n)$, we have $x_n \in \partial f^*(s_n)$ by Proposition C.2.3. The assumption that $s \in \text{int dom } f^*$ means that there exists an $\epsilon > 0$ such that $s + \epsilon\mathbb{B} \subset \text{int dom } f^*$ and f^* is Lipschitz on $s + \epsilon\mathbb{B}$ (Theorem B.3.4). But then $\partial f^*(s + \epsilon\mathbb{B})$ is bounded, and $x_n \in \partial f^*(s_n) \subset \partial f^*(s + \epsilon\mathbb{B})$ for large enough n , contradicting that $\|x_n\| \rightarrow \infty$, and so $s \notin \text{int dom } f^*$ and $\text{int dom } f^* \neq \mathbb{R}^d$.

Now let us assume that $\text{int dom } f^* \neq \mathbb{R}^d$. Let $s \in \text{bd dom } f^*$. Then either $\partial f^*(s) = \emptyset$ or $\partial f^*(s)$ is unbounded (Lemma B.3.21). If $\partial f^*(s) = \emptyset$, take $s_n \rightarrow s$ with $s_n \in \text{relint dom } f^*$, and let $x_n \in \partial f^*(s_n)$. We show that x_n must be unbounded. If x_n is bounded, then by passing to a subsequence if necessary we may assume $x_n \rightarrow x$, and the outer semicontinuity of the subdifferential (Lemma B.3.23) gives $x \in \partial f^*(s)$, contradicting that $\partial f^*(s) = \emptyset$. Thus we must have x_n unbounded, which is thus the desired unbounded sequence. On the other hand, if $\partial f^*(s)$ is unbounded, we can simply take $x_n \in \partial f^*(s)$ with $s \in \partial f(x_n)$ for each n , which is the desired convergent sequence. \square

As a last application of these ideas, in some cases we wish to allow constraints on the functions f to be minimized, returning to the original convex optimization problem (C.0.1) with f a function of Legendre type and C a closed convex set. We then have the following corollary.

Corollary C.2.12. *Let f be of Legendre type (Definition C.3) and $C \subset \mathbb{R}^d$ a closed convex set with $\text{int dom } f \cap C \neq \emptyset$. Define $f_C(x) = f(x) + \mathbf{I}_C(x)$. Then*

(i) f_C^* is continuously differentiable,

(ii) $\text{dom } f_C^* = \mathbb{R}^d$, and

(iii) the constrained tilted minimizers

$$x_u = \operatorname{argmin}_{x \in C} \{\langle u, x \rangle - f(x)\}$$

are unique, continuous in u , belong to $\operatorname{int} \operatorname{dom} f$, and satisfy

$$x_u = \nabla f^*(u - v) \quad \text{and} \quad \nabla f(x_u) = -v$$

for some vector $v \in \mathcal{N}_C(x_u)$.

Proof The function $f_C := f + \mathbf{I}_C$ is closed convex. To show that $\operatorname{dom} f_C^* = \mathbb{R}^d$, we can equivalently show that $(f_C)'_\infty(v) = \infty$ for all non-zero v . Because f is Legendre-type, Lemma C.2.11 guarantees that if $x \in \operatorname{dom} f \cap C$, then

$$(f_C)'_\infty(v) = \lim_{t \uparrow \infty} \frac{f(x + tv) + \mathbf{I}_C(x + tv) - f(x)}{t} \geq \lim_{t \uparrow \infty} \frac{f(x + tv) - f(x)}{t} = f'_\infty(v) = \infty.$$

So $\operatorname{dom} f_C^* = \mathbb{R}^d$, and thus $x_u \operatorname{argmin}_{x \in C} \{f(x) - \langle u, x \rangle\} = \nabla f_C^*(u)$ exists and is unique and continuous, as f is strictly convex.

By the standard subgradient conditions for optimality, the vector x_u is characterized by

$$0 \in \nabla f(x_u) - u + \mathcal{N}_C(x_u),$$

and so $x_u \in \operatorname{int} \operatorname{dom} f$ (as otherwise $\|\nabla f(x_u)\| = +\infty$ by Definition C.3) and

$$x_u = \nabla f^*(u - v)$$

for some vector $v \in \mathcal{N}_C(x_u)$. □

JCD Comment: Now do the particular case that we define $f_C = f + \mathbf{I}_C$ where C is an affine space. Then we should still have $\operatorname{dom} f^* = \mathbb{R}^d$, and ∇f_C^* exists, and should get *some* good dualities. Work it out!

JCD Comment: More smoothness dualities, and write an exercise? Perhaps uniform convexity versions and strict convexity versions.

C.3 Limits at infinity of convex functions and sets

Section C.2 showed the links between conjugate functions and smoothness properties of gradients, making some connections with the existence of minimizers of convex functions. Here, we take this a step further by providing somewhat more abstract conditions for the existence of minimizers of convex functions. Part of this requires a description of different set limits—we wish to understand when we need not take points $\|x\| \rightarrow \infty$ to approach $\inf_x f(x)$ —so we begin there.

C.3.1 Boundedness and closedness of convex sets

The boundedness and closedness of convex sets is central to both the growth properties of convex functions—via their epigraphs—as well as several min-max duality theorems in the coming section. The former will relate to the existence of minimizers of convex functions via so-called recession cones and recession functions associated with convex sets and functions. We develop some of the relevant theory here.

JCD Comment: Draw picture of recession cone

Determining the boundedness and closure of convex sets conveniently often reduces to checking the existence of (infinite) rays remaining within the set. To reduce notational complexity, throughout this section we will assume that $C \subset \mathbb{R}^d$ is a closed convex set, making note when results extend beyond this case. Our starting point is a characterization of various limiting directions that lie within convex sets. For a point $x \in C$, we define the *recession cone*

$$C_\infty(x) := \left\{ v \in \mathbb{R}^d \mid x + tv \in C \text{ for all } t \geq 0 \right\} \quad (\text{C.3.1})$$

to be the set of directions at which C extends off to ∞ . It is immediate that $C_\infty(x)$ is a cone, because $v \in C_\infty(x)$ implies $tv \in C_\infty(x)$ for any $t \geq 0$. As in the case of the recession function (C.3.3), this definition is in fact independent of the choice of x :

Proposition C.3.1. *Let C be closed convex. The set $C_\infty(x)$ is a closed convex cone, independent of x , and*

$$C_\infty(x) = \bigcap_{t>0} \frac{1}{t}(C - x).$$

Proof Let $x_0, x_1 \in C$. For the claim that $C_\infty(x)$ is independent of x , it is enough to show that $C_\infty(x_0) \subset C_\infty(x_1)$. So let $v \in C_\infty(x_0)$ and $t \geq 0$, so that $x_0 + tv \in C$, and for $\epsilon \in [0, 1]$ let

$$y_\epsilon = (1 - \epsilon)x_0 + \epsilon x_1 + tv = (1 - \epsilon) \left(x_0 + \frac{t}{1 - \epsilon} v \right) + \epsilon x_1.$$

Then $y_\epsilon \in C$ for each $\epsilon \in (0, 1)$ because $x_0 + \frac{t}{1 - \epsilon} v \in C$, and $y_1 = \lim_{\epsilon \uparrow 1} y_\epsilon \in C$ as C is closed.

To see the definition of $C_\infty(x)$ in terms of the intersection, note that $v \in C_\infty(x)$ if and only if $tv \in C - x$ for all $t > 0$, that is, $v \in \frac{1}{t}(C - x)$ for all $t > 0$. That $C_\infty(x)$ is thus the intersection of closed convex sets implies it is closed convex. \square

In view of Proposition C.3.1, we can more accurately simply write C_∞ for the recession cone of a convex set C . We can enumerate a few properties of such recession cones.

Lemma C.3.2. *Let C be a closed convex set. Then*

- (i) $C_\infty = \{v \mid C + \mathbb{R}_+ v \subset C\}$ and $C = C + C_\infty$.
- (ii) If $\{C_\alpha\}_{\alpha \in A}$ is a collection of closed convex sets with non-empty intersection, then $(\cap_\alpha C_\alpha)_\infty = \cap_\alpha (C_\alpha)_\infty$.
- (iii) Let u be a unit vector. Then $u \in C_\infty$ if and only if there exists a sequence $x_n \in C$ with $\|x_n\| \rightarrow \infty$ while $x_n / \|x_n\| \rightarrow u$.

(iv) C is bounded if and only if $C_\infty = \{0\}$, that is, it has no directions of recession.

Proof The first claim is immediate from the definition of the recession cone and Proposition C.3.1, so we consider the remainder.

- (ii) We have $v \in (C_\alpha)_\infty$ for each $\alpha \in \mathcal{A}$ if and only if for some $x \in \cap_\alpha C_\alpha$, $x + tv \in C_\alpha$ for each α and $t \geq 0$, that is, $x + tv \in \cap_\alpha C_\alpha$ for all $t \geq 0$ and $v \in (\cap_\alpha C_\alpha)_\infty$.
- (iii) Let $u \in C_\infty$ be a unit vector. Let $x \in C$, and take $x_n = x + nu$ for $n \in \mathbb{N}$. Then clearly $x_n/\|x_n\| \rightarrow u$, while $x_n \in C$ for each n . Conversely, assume that $x_n/\|x_n\| \rightarrow u$ while $\|x_n\| \rightarrow \infty$ and $x_n \in C$. Let $t \geq 0$ and N be large enough that $\|x_n\| \geq t$ for $n \geq N$. Then

$$x + tu = \lim_n \underbrace{\left(1 - \frac{t}{\|x_n\|}\right)x + \frac{t}{\|x_n\|}x_n}_{\text{in } C \text{ for } n \geq N}.$$

That is, $x + tu$ is the limit of points in C , so that the closedness of C implies $x + tu \in C$.

- (iv) We certainly have $C_\infty = \{0\}$ when C is bounded. If C is unbounded, then we may take a sequence $x_n \in C$ with $\|x_n\| \rightarrow \infty$, and by considering convergent subsequences, may assume w.l.o.g. that $x_n/\|x_n\| \rightarrow u$ for some $u \in \mathbb{S}^{d-1}$. Then $C_\infty \supset \{u\}$ is non-empty by part (iii). □

JCD Comment: Draw a figure of Theorem C.3.3

Motivated by the conjugate duality relationships that Theorem C.2.1 presents—so that $f = f^{**}$ if and only if f is closed—we will find it useful to consider closedness properties of convex sets C under linear mappings. Recall that for $A \in \mathbb{R}^{n \times d}$, we write $AC = \{Ax \mid x \in C\}$, and let $\text{null}(A) = \{z \mid Az = 0\}$ be its nullspace. We have the following theorem on closedness.

Theorem C.3.3. *Let $C \subset \mathbb{R}^d$ be a closed convex set and $A \in \mathbb{R}^{m \times d}$. If $C_\infty \cap \text{null}(A)$ is a linear subspace, then AC is closed.*

Proof Define the set

$$L := C_\infty \cap \text{null}(A),$$

which is a linear subspace and evidently satisfies $L = (-C_\infty) \cap C_\infty \cap \text{null}(A)$ as well. Then $L^\perp = \{v \mid \langle v, z \rangle = 0 \text{ for } z \in L\}$ satisfies

$$C = (L^\perp \cap C) + L,$$

because $L \subset C_\infty \cap (-C_\infty)$ and $C = C + L$ (recall Lemma C.3.2, part (i)). Therefore, $AC = A(L^\perp \cap C)$ as $AL = \{0\}$. Now let $y \in \text{cl}(AC)$. For $\epsilon > 0$, define the sets

$$B_\epsilon := \{x \mid \|Ax - y\| \leq \epsilon\} \quad \text{and} \quad C_\epsilon := C \cap L^\perp \cap B_\epsilon = \left\{x \in C \cap L^\perp \mid \|Ax - y\| \leq \epsilon\right\}.$$

Because $AC = A(L^\perp \cap C)$, the sets C_ϵ are non-empty, closed, and nested. If we can show that C_ϵ is bounded, then it will be compact, and the finite intersection property will imply that $\cap_{\epsilon > 0} C_\epsilon \neq \emptyset$, and so there will be some $x \in C$ for which $Ax = y$.

To show the boundedness, we use Lemma C.3.2. Because $(B_\epsilon)_\infty = \text{null}(A)$, part (ii) of Lemma C.3.2 gives

$$(C_\epsilon)_\infty = C_\infty \cap L_\infty^\perp \cap (B_\epsilon)_\infty = C_\infty \cap L^\perp \cap \text{null}(A) = L^\perp \cap L = \{\mathbf{0}\},$$

so that part (iv) of the same lemma implies C_ϵ is bounded and hence compact, as desired. \square

As a representative corollary of Theorem C.3.3, let us revisit the problem of partial minimization, as in Proposition B.3.11. Let $F : \mathbb{R}^d \times \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ be convex, and for $x \in \mathbb{R}^d, y \in \mathbb{R}^n$ define the partial minimization function

$$f(x) := \inf_y F(x, y).$$

We have seen already that f is convex—we wish to know whether f is closed convex. A sufficient condition here is that F be closed convex and $F(x, \cdot)$ be coercive for some x .

Corollary C.3.4. *Let F be closed convex. If there exists x_0, t_0 for which the set*

$$\{y \mid F(x_0, y) \leq t_0\}$$

is bounded, then $f(x) = \inf_y F(x, y)$ is closed convex.

Proof Consider the epigraphs

$$\text{epi } f = \{(x, t) \mid f(x) \leq t\} = \left\{ (x, t) \mid \inf_y F(x, y) \leq t \right\}$$

and

$$\text{epi } F = \{(x, y, t) \mid F(x, y) \leq t\}.$$

Define the projection $\pi : \mathbb{R}^d \times \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}^d \times \mathbb{R}$ by $\pi(x, y, t) = (x, t)$ (which is evidently a linear operator). Then we claim that

$$\pi(\text{epi } F) \subset \text{epi } f \subset \text{cl } \pi(\text{epi } F). \quad (\text{C.3.2})$$

Indeed, if $(x, y, t) \in \text{epi } F$, then $(x, t) \in \text{epi } f$ trivially, showing the first inclusion. For the second, note that if $(x, t) \in \text{epi } f$, then there must be a sequence y_n $\liminf_n F(x, y_n) \leq t$. Then for any $\epsilon > 0$, $(x, y_n, t + \epsilon) \in \text{epi } F$ for all large n , and $(x, t + \epsilon) \in \pi(\text{epi } F)$. So $(x, t) \in \text{cl } \pi(\text{epi } F)$.

In view of the containments (C.3.2), to show that $\text{epi } f$ is closed it is enough to show that $\pi(\text{epi } F)$ is closed. By Theorem C.3.3, for this it is sufficient to show that $(\text{epi } F)_\infty \cap \text{null}(\pi)$ is a linear subspace, and using that $\text{null}(\pi) = \{(\mathbf{0}, y, 0) \mid y \in \mathbb{R}^n\}$, this set is

$$\begin{aligned} (\text{epi } F)_\infty \cap \text{null}(\pi) &= \{(\mathbf{0}, v_y, 0) \mid (0, v_y, 0) \in (\text{epi } F)_\infty\} \\ &= \{(\mathbf{0}, v_y, 0) \mid \text{there exist } x, y, r \text{ s.t. } (x, y + tv_y, r) \in \text{epi } F \text{ for all } t \geq 0\}. \end{aligned}$$

But as is now familiar from our treatment of recession functions and recession cones, we know that in the final set, the choices $(x, y, r) \in \text{epi } F$ are arbitrary (Proposition C.3.1), and in particular, we may take $x = x_0$ for the x_0 making the set $\{y \mid F(x_0, y) \leq t_0\}$ compact. But Lemma C.3.2, part (iv) guarantees that only $v_y = \mathbf{0}$ satisfies $(x_0, y + tv_y, r) \in \text{epi } F$ for all $t \geq 0$. \square

C.3.2 Asymptotic growth and existence of minimizers

We can use the identification between the domains of ∂f and the images of ∂f^* to give a few additional characterizations of the domains of convex functions and their conjugates; the domain of f^* is intimately tied with the growth properties of f , and conversely by the relationship $f = f^{**}$ when f is closed convex. As one example of how we can make this identification, note that if f^* is defined everywhere, that is, $\text{dom } f^* = \mathbb{R}^d$, then similarly $\text{dom } \partial f^* = \mathbb{R}^d$, and so in particular the (sub)gradients of f must cover all of \mathbb{R}^d . Even more, as we shall see, this implies certain growth conditions on f .

To make this more rigorous, we require functions capturing the asymptotic growth of f , analogizing the recession cones (C.3.1). To that end, we present the following proposition, which has the benefit of defining the *recession function* (essentially, an asymptotic derivative) of f .

Proposition C.3.5. *Let f be a closed convex function and f^* is convex conjugate. Then for any $x \in \text{dom } f$, we may define*

$$f'_\infty(v) := \sup_{t>0} \frac{f(x+tv) - f(x)}{t} = \lim_{t \rightarrow \infty} \frac{f(x+tv) - f(x)}{t} \quad (\text{C.3.3})$$

independently of x , and moreover,

$$f'_\infty(v) = \sigma_{\text{dom } f^*}(v)$$

where $\sigma_{\text{dom } f^*}$ is the support function (B.2.1) of $\text{dom } f^*$.

Proof That for any fixed $x \in \text{dom } f$ the limit exists and is equal to the supremum follows because of the criterion of increasing slopes (B.3.4), making the equality with the supremum immediate. That $f'_\infty(v)$ is independent of x will follow once we show the second equality claimed in the proposition, to which we now turn.

Recall that

$$\text{dom } f^* = \left\{ s \mid \sup_x \{ \langle s, x \rangle - f(x) \} < \infty \right\} \quad \text{and} \quad f(x) = \sup_s \{ \langle s, x \rangle - f^*(s) \}$$

by conjugate duality, as f is closed convex (Theorem C.2.1). Fix $x \in \text{dom } f$. Then for any $s \in \text{dom } f^*$, we evidently have

$$\frac{f(x+tv) - f(x)}{t} \geq \frac{\langle s, x+tv \rangle - f^*(s) - f(x)}{t} \rightarrow \langle s, v \rangle$$

as $t \uparrow \infty$. Taking a supremum over $s \in \text{dom } f^*$ gives that $f'_\infty(v) \geq \sigma_{\text{dom } f^*}(v)$. For the opposite direction, note that

$$\begin{aligned} \frac{f(x+tv) - f(x)}{t} &= \frac{1}{t} \left[\sup_{s \in \text{dom } f^*} \{ \langle s, x+tv \rangle - f^*(s) \} - \sup_{s \in \text{dom } f^*} \{ \langle s, x \rangle - f^*(s) \} \right] \\ &\leq \frac{1}{t} \sup_{s \in \text{dom } f^*} \{ \langle s, x+tv \rangle - f^*(s) - (\langle s, x \rangle - f^*(s)) \} = \frac{1}{t} \sup_{s \in \text{dom } f^*} t \langle s, v \rangle. \end{aligned}$$

Thus $f'_\infty(v) \leq \sigma_{\text{dom } f^*}(v)$, and we have the result. \square

JCD Comment: See Theorem 13.3 of Rockafellar.

It is particularly interesting to understand the conditions under which $\text{dom } f^* = \mathbb{R}^d$, that is, f^* is finite everywhere, and relatedly, under which the function $x \mapsto f(x) - \langle s, x \rangle$ has a minimizer. Recall that $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ is *coercive* if $f(x) \rightarrow \infty$ whenever $\|x\| \rightarrow \infty$. For convex functions, we may characterize coercivity by the recession function of f .

Proposition C.3.6. *Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be closed convex. Then f is coercive if and only if $f'_\infty(v) > 0$ for all $v \neq 0$, that is, f is coercive on each line.*

Proof If f is coercive, then for any $v \neq 0$, $f(x + tv) \rightarrow \infty$ as $t \rightarrow \infty$, so that for any $x \in \text{dom } f$, $f(x + tv) - f(x) > 0$ for large enough t . In particular,

$$f'_\infty(v) = \sup_{t>0} \frac{f(x + tv) - f(x)}{t} > 0$$

by Proposition C.3.5. Conversely, assume that $f'_\infty(v) > 0$ for all $v \neq 0$. Then if f is not coercive, it must have an unbounded level set $S_\alpha := \{x \mid f(x) \leq \alpha\}$. As S_α is a convex set, if it is unbounded, then there exists a unit vector $v \in (S_\alpha)_\infty$, the recession cone of S_α by Lemma C.3.2. But then $x + tv \in S_\alpha$ for all $t \geq 0$ implies that $f(x + tv) \leq \alpha$ for all $t \geq 0$, a contradiction to $f'_\infty(v) > 0$. \square

We call f *super-coercive* if $f(x)/\|x\| \rightarrow \infty$ whenever $\|x\| \rightarrow \infty$, so that f grows more than linearly. These concepts are central to the existence of minimizers. A priori, any function with compact domain is super-coercive, because $f(x) = +\infty$ for $x \notin \text{dom } f$. The recession function f'_∞ associated with f , as expression (C.3.3) defines, characterizes these functions as well. Particularly important are those f satisfying

$$f'_\infty(v) = +\infty \text{ for all } v \neq 0,$$

a class Rockafellar [162] calls *copositive* functions, as these exhibit superlinear growth on all rays toward infinity. We can relate this condition to the domains of f^* as well: using Proposition B.2.7, Propositions C.3.5 and C.3.6 give

Corollary C.3.7. *Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be closed convex. Then $s \in \text{dom } f^*$ if and only if $\langle s, v \rangle \leq f'_\infty(v)$ for all $v \neq 0$, and $s \in \text{int dom } f^*$ if and only if $\langle s, v \rangle < f'_\infty(v)$ for all $v \neq 0$. In particular,*

(i) *If f is coercive, or, equivalently, if $f'_\infty(v) > 0$ for all $v \neq 0$, then $0 \in \text{int dom } f^*$ and f has a minimizer.*

(ii) *We have $f'_\infty(v) = +\infty$ for all $v \neq 0$ if and only if*

$$\text{dom } f^* = \mathbb{R}^d.$$

A sufficient condition for this is that f be super-coercive.

Proof Combine Propositions B.2.7, C.3.5, and C.3.6: for part (i), note that if $f'_\infty(v) > 0$ for all $v \neq 0$, then $0 \in \text{int dom } f^*$, and so f^* has a non-trivial subdifferential $\partial f^*(0)$ at 0; letting $x \in \partial f^*(0)$ we have $x \in \text{argmin } f$. Part (ii) is similarly immediate. \square

As an example consequence, if f is closed convex and $c = \inf_{\|v\|=1} f'_\infty(v) > 0$ is the minimal value of the recession function, then the tilted function $f(\cdot) - \langle s, \cdot \rangle$ has a minimizer whenever $\|s\| < c$.

C.4 Saddle point theorems and min-max duality

Our final set of results centers around min-max duality theorems, central concepts in convex analysis that, for us, are important in the development of proper losses (Chapter 14) and minimax games (Chapter 19). In the main, we will consider conditions under which, for a function $L : \mathbb{R}^d \times \mathbb{R}^k \rightarrow \mathbb{R}$, we have

$$\inf_{x \in X} \sup_{y \in Y} L(x, y) = \sup_{y \in Y} \inf_{x \in X} L(x, y). \quad (\text{C.4.1})$$

We begin with a few straightforward remarks giving sufficient conditions to swap the min and max (i.e., the infimum and supremum) in the equality (C.4.1), then develop more general theory that shows when this is indeed possible. When equality (C.4.1) holds, we shall say that *strong duality holds* in the saddle point.

The starting point is the weak min-max inequality, which holds without conditions on L or the sets X, Y .

Proposition C.4.1 (The weak min-max inequality). *For any sets X, Y and function L ,*

$$\sup_{y \in Y} \inf_{x \in X} L(x, y) \leq \inf_{x \in X} \sup_{y \in Y} L(x, y). \quad (\text{C.4.2})$$

Proof Fix $y_0 \in Y$. Then for any $x_0 \in X$, $\inf_{x \in X} L(x, y_0) \leq L(x_0, y_0) \leq \sup_{y \in Y} L(x_0, y)$. Taking a supremum on the left side implies $\sup_{y \in Y} \inf_{x \in X} L(x, y) \leq \sup_{y \in Y} L(x_0, y)$ for any $x_0 \in X$. We may now take an infimum on the right. \square

The simplest condition guaranteeing equality (C.4.1) is the existence of a saddle point $(x^*, y^*) \in X \times Y$, meaning a point satisfying

$$\sup_{y \in Y} L(x^*, y) \leq L(x^*, y^*) \leq \inf_{x \in X} L(x, y^*). \quad (\text{C.4.3})$$

One of the main concerns of duality theory and min-max games is when such points exist, as they guarantee the equality (C.4.1) and even achieve the values of the game:

Proposition C.4.2. *A pair $(x^*, y^*) \in X \times Y$ is a saddle point for L on $X \times Y$ if and only if all of the following are true:*

- i. *the strong duality (C.4.1) holds*
- ii. *x^* minimizes $\sup_{y \in Y} L(x, y)$ over $x \in X$*
- iii. *y^* maximizes $\inf_{x \in X} L(x, y)$ over $y \in Y$*

Under these conditions, the value of the game is

$$L(x^*, y^*) = \inf_{x \in X} L(x, y^*) = \sup_{y \in Y} L(x^*, y) = \inf_{x \in X} \sup_{y \in Y} L(x, y) = \sup_{y \in Y} \inf_{x \in X} L(x, y).$$

Proof Given a saddle point (x^*, y^*) , we have

$$\inf_{x \in X} \sup_{y \in Y} L(x, y) \leq \sup_{y \in Y} L(x^*, y) \leq L(x^*, y^*) \leq \inf_{x \in X} L(x, y^*) \leq \sup_{y \in Y} \inf_{x \in X} L(x, y),$$

so that in view of the weak min-max inequality (C.4.2), each inequality must actually be an equality, and equality (C.4.1) holds as well. Thus, $\inf_{x \in X} L(x, y^*) = \sup_{y \in Y} \inf_{x \in X} L(x, y)$, and y^* minimizes $\inf_{x \in X} L(x, y)$ over $y \in Y$, and similarly for x^* .

Now suppose conditions i–iii hold above. Then

$$\inf_{x \in X} \sup_{y \in Y} L(x, y) = \sup_{y \in Y} L(x^*, y) \geq L(x^*, y^*) \geq \inf_{x \in X} L(x, y^*) = \sup_{y \in Y} \inf_{x \in X} L(x, y).$$

By assumption, the left and right quantities are equal, so that (x^*, y^*) is a saddle point. \square

C.4.1 Saddle points and convex conjugates

We now proceed to more general considerations under which equality (C.4.1) holds. For the remainder of the section, we shall assume that L is a closed convex-concave function on $X \times Y$, meaning that for each $y \in Y$,

$$x \mapsto L(x, y) + \mathbf{I}_X(x) = \begin{cases} L(x, y) & \text{if } x \in X \\ +\infty & \text{otherwise} \end{cases}$$

is a closed convex function, while for each $x \in X$,

$$y \mapsto -L(x, y) + \mathbf{I}_Y(y) = \begin{cases} -L(x, y) & \text{if } y \in Y \\ +\infty & \text{otherwise} \end{cases}$$

is closed convex. We will extend our arguments with tilted minimizers and convex conjugates, so that the key object we consider will be the parameterized function

$$F(x, u) := \sup_{y \in Y} \{L(x, y) + \langle u, y \rangle\}.$$

Evidently, F is a closed convex function whenever L is closed convex-concave, as it is the supremum of closed convex functions. We also define the parameterized (primal) value function

$$p(u) := \inf_{x \in X} F(x, u), \tag{C.4.4}$$

which is a convex function in u (as it is a partial minimization), and satisfies

$$p(0) = \inf_{x \in X} \sup_{y \in Y} L(x, y).$$

Our first step is to derive the conjugate of this primal value function.

Lemma C.4.3. *Let $p^*(v) = \sup_u \{\langle u, v \rangle - p(u)\}$. Then*

$$p^*(v) = \begin{cases} -\inf_{x \in X} L(x, v) & \text{if } v \in Y \\ +\infty & \text{otherwise.} \end{cases}$$

Proof We expand

$$\begin{aligned} p^*(v) &= \sup_u \left\{ \langle u, v \rangle - \inf_{x \in X} \sup_{y \in Y} \{L(x, y) + \langle u, y \rangle\} \right\} \\ &= \sup_u \sup_{x \in X} \inf_{y \in Y} \{ \langle u, v - y \rangle - L(x, y) \} = \sup_{x \in X} \sup_u \inf_{y \in Y} \{ \langle u, v - y \rangle - L(x, y) \}. \end{aligned}$$

By assumption, for each $x \in X$ the epigraph

$$E_x := \{(y, t) \in Y \times \mathbb{R} \mid -L(x, y) \leq t\} \quad (\text{C.4.5})$$

is a closed convex set.

Let us first assume $v \in Y$. Then there exists a non-vertical supporting hyperplane $(w, 1)$ to E_x at the point $(v, -L(x, v))$, that is, a $w \in \mathbb{R}^k$ satisfying

$$\langle w, v \rangle - L(x, v) = \inf_{(y, t) \in E_x} \{ \langle w, y \rangle - t \} = \inf_{y \in Y} \{ \langle w, y \rangle - L(x, y) \}.$$

Of course, for $v \in Y$, we may always choose $y = v$ in the infimum and always have

$$\inf_{y \in Y} \{ \langle u, y - v \rangle - L(x, y) \} \leq \langle u, v - v \rangle - L(x, v) = -L(x, v).$$

Equality obtains there for $w = u$, so

$$\sup_u \inf_{y \in Y} \{ \langle u, v - y \rangle - L(x, y) \} = \inf_{y \in Y} \{ \langle w, v - y \rangle - L(x, y) \} = -L(x, v),$$

which gives the first claim of the equality defining $p^*(v)$.

When $v \notin Y$, then evidently for all $x \in X$ and $t \in \mathbb{R}$, we have $(v, t) \notin E_x$, recall the epigraph (C.4.5). So for any $t \in \mathbb{R}$ and $x \in X$, there exists a non-vertical hyperplane strictly separating (v, t) from E_x (Corollary B.1.12). That is, we have a $w \in \mathbb{R}^k$ such that

$$t + \langle w, v \rangle < \inf_{y \in Y} \{ \langle w, y \rangle - L(x, y) \} \quad \text{i.e.} \quad t < \inf_{y \in Y} \{ \langle w, y - v \rangle - L(x, y) \}.$$

Because $t < \infty$ was arbitrary, we must have $\sup_u \inf_{y \in Y} \{ \langle u, v - y \rangle - L(x, y) \} = +\infty$. □

Given Lemma C.4.3, we can almost immediately see that the biconjugate of p will determine whether we have the saddle point equality (C.4.1). We state this as a proposition.

Proposition C.4.4. *Let p be the primal value function (C.4.4). Then*

$$p^{**}(0) = \sup_{y \in Y} \inf_{x \in X} L(x, y).$$

The saddle point equality (C.4.1) holds if p is lower semicontinuous at 0 with either $p(0) < \infty$ or $p(0) = \infty$ and $p(u) > -\infty$ for all u .

Proof Taking the biconjugate, we have

$$p^{**}(u) = \sup_y \{ \langle y, u \rangle - p^*(y) \} = \sup_{y \in Y} \left\{ \langle y, u \rangle + \inf_{x \in X} L(x, y) \right\}$$

by Lemma C.4.3. This shows the first equality of the proposition. For the second claim, apply Corollary C.2.2, which shows that $p^{**}(0) = p(0)$ if and only if p is lower semicontinuous at 0. □

C.4.2 Min-max duality and the existence of saddle points

Recalling the definition (C.4.4) of the value function via $F(x, u) = \sup_{y \in Y} \{L(x, y) + \langle u, y \rangle\}$, we are led back to investigate when exactly partial minimization of convex functions yields a closed convex function, as $p(u) = \inf_{x \in X} F(x, u)$. We can provide a few consequences of Proposition C.4.4 that guarantee strong duality holds for the saddle point problem (C.4.1).

Theorem C.4.5. *Let X, Y be convex sets and let L be closed convex-concave on $X \times Y$. Assume that*

$$f(x) := \begin{cases} \sup_{y \in Y} L(x, y) & \text{if } x \in X \\ +\infty & \text{otherwise} \end{cases}$$

satisfies $\inf_x f(x) < \infty$ and has compact level sets. Then the value function (C.4.4) is closed convex, strong duality holds for the saddle point problem (C.4.1), and there exists a nonempty compact set $X^ \subset X$ minimizing f .*

Proof By Corollary C.3.4, for $F(x, u) = \sup_{y \in Y} \{L(x, y) + \langle u, y \rangle\} + \mathbf{I}_X(x)$, it is enough to show that there exists u_0 such that $x \mapsto F(x, u_0)$ has compact level sets. Taking $u = 0$, we have $F(x, 0) = f(x)$, and so the level sets $\{x \mid F(x, 0) \leq t\}$ are compact by assumption. So $p(u)$ is lower semicontinuous in u , and $p^{**}(u) = p(u)$ for all u . Apply Proposition C.4.4 to obtain strong duality.

For the existence of a set of minimizers, note that $f(x)$ is closed convex and has compact level sets, so its minimizers are attained. \square

In brief, under the conditions of the theorem, if we let $x^* \in X^*$, we have for any $y \in Y$ that

$$\sup_{y \in Y} L(x^*, y) = \inf_{x \in X} \sup_{y \in Y} L(x, y) = \sup_{y \in Y} \inf_{x \in X} L(x, y) \leq \sup_{y \in Y} L(x^*, y),$$

and equality must hold in the final inequality. That is, we have a sort of “partial” saddle point: there exists x^* achieving the value of the game that L implies, so

$$\sup_{y \in Y} L(x^*, y) = \sup_{y \in Y} \inf_{x \in X} L(x, y).$$

We also have the following corollary.

Corollary C.4.6. *Let X, Y be convex sets, and let X be compact. Assume L is closed convex-concave on $X \times Y$. Then if there exists $x \in X$ such that $\sup_{y \in Y} L(x, y) < \infty$,*

$$\inf_{x \in X} \sup_{y \in Y} L(x, y) = \sup_{y \in Y} \inf_{x \in X} L(x, y)$$

and there exists a non-empty convex compact set X^ minimizing $\sup_{y \in Y} L(x, y)$ over $x \in X$.*

An elaborated version of these argument guarantees the existence of saddle points, so that strong duality (C.4.1) obtains.

Theorem C.4.7. *Let X, Y be convex sets and L be closed convex-concave and finite on $X \times Y$. Then the set $X^* \times Y^*$ of saddle points for L on $X \times Y$ is convex, compact, and non-empty if any of the following conditions hold:*

- i. The functions $f(x) := \sup_{y \in Y} L(x, y) + \mathbf{I}_X(x)$ and $g(y) := \sup_{x \in X} -L(x, y) + \mathbf{I}_Y(y)$ are coercive.*

ii. The sets X and Y are compact.

Proof Theorem C.4.5 guarantees that $\inf_{x \in X} \sup_{y \in Y} L(x, y) = \sup_{y \in Y} \inf_{x \in X} \sup_{y \in Y} L(x, y)$ under any of the conditions. It also guarantees the existence of a compact set of minimizers $X^* \subset X$ such that

$$\sup_{y \in Y} L(x^*, y) = \inf_{x \in X} \sup_{y \in Y} L(x, y)$$

for $x^* \in X^*$. Flipping the roles of x and y , we see similarly that there exists a compact convex set $Y^* \subset Y$ such that

$$\sup_{x \in X} -L(x, y^*) = \inf_{y \in Y} \sup_{x \in X} -L(x, y),$$

that is, $\inf_{x \in X} L(x, y^*) = \sup_{y \in Y} \inf_{x \in X} L(x, y)$. Proposition C.4.2 gives the theorem. \square

Further reading

There are a variety of references on the topic, beginning with the foundational book by Rockafellar [162], which initiated the study of convex functions and optimization in earnest. Since then, a variety of authors have written (perhaps more easily approachable) books on convex functions, optimization, and their related calculus. Hiriart-Urruty and Lemaréchal [111] have written two volumes explaining in great detail finite-dimensional convex analysis, and provide a treatment of some first-order algorithms for solving convex problems. Borwein and Lewis [36] and Luenberger [141] give general treatments that include infinite-dimensional convex analysis, and Bertsekas [29] gives a variety of theoretical results on duality and optimization theory.

There are, of course, books that combine theoretical treatment with questions of convex modeling and procedures for solving convex optimization problems (problems for which the objective and constraint sets are all convex). Boyd and Vandenberghe [38] gives a very readable treatment for those who wish to use convex optimization techniques and modeling, as well as the basic results in convex analytic background and duality theory. Ben-Tal and Nemirovski [22], as well as Nemirovski's various lecture notes, give a theory of the tractability of computing solutions to convex optimization problems as well as methods for solving them.

C.5 Exercises

Exercise C.1: Show that the alternative increasing slopes condition (B.3.5) is equivalent to convexity of f .

Exercise C.2: Prove Proposition C.1.7. *Hint:* See the proof of Proposition C.1.5.

Exercise C.3: Prove Proposition C.1.8. *Hint:* See the proof of Proposition C.1.6.

Exercise C.4: Do the uniform convexity version of Proposition C.1.5.

Exercise C.5: Do the uniform convexity version of Proposition C.1.6.

Bibliography

- [1] F. Abramovich, Y. Benjamini, D. L. Donoho, and I. M. Johnstone. Adapting to unknown sparsity by controlling the false discovery rate. *Annals of Statistics*, 34(2):584–653, 2006.
- [2] D. Achlioptas. Database-friendly random projections: Johnson-Lindenstrauss with binary coins. *Journal of Computer and System Sciences*, 66:671–687, 2003.
- [3] S. Agrawal and N. Goyal. Analysis of Thompson sampling for the multi-armed bandit problem. In *Proceedings of the Twenty Fifth Annual Conference on Computational Learning Theory*, 2012.
- [4] R. Ahlswede, P. Gács, and J. Körner. Bounds on conditional probabilities with applications in multi-user communication. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*, 34(2):157–177, 1976.
- [5] N. Ailon and B. Chazelle. The fast Johnson-Lindenstrauss transform and approximate nearest neighbors. *SIAM Journal on Computing*, 39(1):302–322, 2009.
- [6] S. M. Ali and S. D. Silvey. A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society, Series B*, 28:131–142, 1966.
- [7] A. Andoni, M. Datar, N. Immorlica, P. Indyk, and V. Mirrokni. Locality-sensitive hashing using stable distributions. In T. Darrell, P. Indyk, and G. Shakhnarovich, editors, *Nearest Neighbor Methods in Learning and Vision: Theory and Practice*. MIT Press, 2006.
- [8] E. Arias-Castro, E. Candés, and M. Davenport. On the fundamental limits of adaptive sensing. *IEEE Transactions on Information Theory*, 59(1):472–481, 2013.
- [9] S. Artstein, K. Ball, F. Barthe, and A. Naor. Solution of Shannon’s problem on the monotonicity of entropy. *Journal of the American Mathematical Society*, 17(4):975–982, 2004.
- [10] P. Assouad. Deux remarques sur l’estimation. *Comptes Rendus des Séances de l’Académie des Sciences, Série I*, 296(23):1021–1024, 1983.
- [11] J.-Y. Audibert and S. Bubeck. Regret bounds and minimax policies under partial monitoring. In *Journal of Machine Learning Research*, pages 2635–2686, 2010.
- [12] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2-3):235–256, 2002.
- [13] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.

- [14] B. Balle, G. Barthe, and M. Gaboardi. Privacy amplification by subsampling: Tight analyses via couplings and divergences. In *Advances in Neural Information Processing Systems 31*, pages 6277–6287, 2018.
- [15] Z. Bar-Yossef, T. S. Jayram, R. Kumar, and D. Sivakumar. An information statistics approach to data stream and communication complexity. *Journal of Computer and System Sciences*, 68(4):702–732, 2004.
- [16] A. Barron. Entropy and the central limit theorem. *Annals of Probability*, 14(1):336–342, 1986.
- [17] A. R. Barron. Complexity regularization with application to artificial neural networks. In *Nonparametric Functional Estimation and Related Topics*, pages 561–576. Kluwer Academic, 1991.
- [18] A. R. Barron and T. M. Cover. Minimum complexity density estimation. *IEEE Transactions on Information Theory*, 37:1034–1054, 1991.
- [19] P. L. Bartlett, M. I. Jordan, and J. McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101:138–156, 2006.
- [20] R. Bassily, A. Smith, T. Steinke, and J. Ullman. More general queries and less generalization error in adaptive data analysis. *arXiv:1503.04843 [cs.LG]*, 2015.
- [21] R. Bassily, K. Nissim, A. Smith, T. Steinke, U. Stemmer, and J. Ullman. Algorithmic stability for adaptive data analysis. In *Proceedings of the Forty-Eighth Annual ACM Symposium on the Theory of Computing*, pages 1046–1059, 2016.
- [22] A. Ben-Tal and A. Nemirovski. *Lectures on Modern Convex Optimization*. SIAM, 2001.
- [23] Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society, Series B*, 57(1):289–300, 1995.
- [24] D. Berend and A. Kontorovich. On the concentration of the missing mass. *Electronic Communications in Probability*, 18:1–7, 2018.
- [25] J. O. Berger. The case for objective Bayesian analysis. *Bayesian Analysis*, 1(3):385–402, 2006.
- [26] J. O. Berger, J. Bernardo, and D. Sun. The formal definition of reference priors. *Annals of Statistics*, 37(2):905–938, 2009.
- [27] R. Berk, L. Brown, A. Buja, K. Zhang, and L. Zhao. Valid post-selection inference. *Annals of Statistics*, 41(2):802–837, 2013.
- [28] J. M. Bernardo. Reference analysis. In D. Day and C. R. Rao, editors, *Bayesian Thinking, Modeling and Computation*, volume 25 of *Handbook of Statistics*, chapter 2, pages 17–90. Elsevier, 2005.
- [29] D. P. Bertsekas. *Convex Optimization Theory*. Athena Scientific, 2009.
- [30] P. Billingsley. *Convergence of Probability Measures*. Wiley, Second edition, 1999.

- [31] L. Birgé. Approximation dans les espaces métriques et théorie de l'estimation. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*, 65:181–238, 1983.
- [32] L. Birgé. A new lower bound for multiple hypothesis testing. *IEEE Transactions on Information Theory*, 51(4):1611–1614, 2005.
- [33] L. Birgé and P. Massart. Estimation of integral functionals of a density. *Annals of Statistics*, 23(1):11–29, 1995.
- [34] J. Blasiok, P. Gopalan, L. Hu, and P. Nakkiran. A unifying theory of distance from calibration. In *Proceedings of the Fifty-Fifth Annual ACM Symposium on the Theory of Computing*, 2023. URL <https://arxiv.org/abs/2211.16886>.
- [35] J. Blasiok, P. Gopalan, L. Hu, and P. Nakkiran. When does optimizing a proper loss yield calibration? In *Advances in Neural Information Processing Systems 36*, 2023. URL <https://arxiv.org/abs/2305.18764>.
- [36] J. Borwein and A. Lewis. *Convex Analysis and Nonlinear Optimization*. Springer, 2006.
- [37] S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: a Nonasymptotic Theory of Independence*. Oxford University Press, 2013.
- [38] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [39] M. Braverman, A. Garg, T. Ma, H. L. Nguyen, and D. P. Woodruff. Communication lower bounds for statistical estimation problems via a distributed data processing inequality. In *Proceedings of the Forty-Eighth Annual ACM Symposium on the Theory of Computing*, 2016. URL <https://arxiv.org/abs/1506.07216>.
- [40] G. Brown, M. Bun, V. Feldman, A. Smith, and K. Talwar. When is memorization of irrelevant training data necessary for high-accuracy learning? In *Proceedings of the Fifty-Third Annual ACM Symposium on the Theory of Computing*, pages 123–132, 2021.
- [41] L. D. Brown. *Fundamentals of Statistical Exponential Families*. Institute of Mathematical Statistics, Hayward, California, 1986.
- [42] S. Bubeck and N. Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.
- [43] V. Buldygin and Y. Kozachenko. *Metric Characterization of Random Variables and Random Processes*, volume 188 of *Translations of Mathematical Monographs*. American Mathematical Society, 2000.
- [44] T. Cai and M. Low. Testing composite hypotheses, Hermite polynomials and optimal estimation of a nonsmooth functional. *Annals of Statistics*, 39(2):1012–1041, 2011.
- [45] T. T. Cai, Y. Wang, and L. Zhang. The cost of privacy: optimal rates of convergence for parameter estimation with differential privacy. *Annals of Statistics*, 49(5):2825–2850, 2021.
- [46] T. T. Cai, Y. Wang, and L. Zhang. Score attack: A lower bound technique for optimal differentially private learning. *arXiv:2303.07152 [math.ST]*, 2023.

- [47] E. J. Candès and M. A. Davenport. How well can we estimate a sparse vector. *Applied and Computational Harmonic Analysis*, 34(2):317–323, 2013.
- [48] O. Catoni. *PAC-Bayesian Supervised Classification: The Thermodynamics of Statistical Learning*, volume 56 of *IMS Lecture Notes and Monographs*. Institute of Mathematical Statistics, Beachwood, Ohio, USA, 2007. URL <https://arxiv.org/abs/0712.0248>.
- [49] O. Catoni and I. Giulini. Dimension-free PAC-bayesian bounds for matrices, vectors, and linear least squares regression. *arXiv:1712.02747 [math.ST]*, 2017.
- [50] N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, 2006.
- [51] O. Chapelle and L. Li. An empirical evaluation of Thompson sampling. In *Advances in Neural Information Processing Systems 24*, 2011.
- [52] B. S. Clarke and A. R. Barron. Information-theoretic asymptotics of bayes methods. *IEEE Transactions on Information Theory*, 36(3):453–471, 1990.
- [53] B. S. Clarke and A. R. Barron. Jeffreys’ prior is asymptotically least favorable under entropy risk. *Journal of Statistical Planning and Inference*, 41:37–60, 1994.
- [54] J. E. Cohen, Y. Iwasa, G. Rautu, M. B. Ruskai, E. Seneta, and G. Zbaganu. Relative entropy under mappings by stochastic matrices. *Linear Algebra and its Applications*, 179:211–235, 1993.
- [55] C. Cortes and V. Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, September 1995.
- [56] J. Couzin. Whole-genome data not anonymous, challenging assumptions. *Science*, 321(5894):1278, 2008.
- [57] T. M. Cover and J. A. Thomas. *Elements of Information Theory, Second Edition*. Wiley, 2006.
- [58] K. Crammer and Y. Singer. On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of Machine Learning Research*, 2:265–292, 2001.
- [59] I. Csiszár. Information-type measures of difference of probability distributions and indirect observation. *Studia Scientifica Mathematica Hungary*, 2:299–318, 1967.
- [60] I. Csiszár. A class of measures of informativity of observation channels. *Periodica Mathematica Hungarica*, 2(1–4):191–213, 1972.
- [61] I. Csiszár and J. Körner. *Information Theory: Coding Theorems for Discrete Memoryless Systems*. Cambridge University Press, second edition, 2011.
- [62] T. Dalenius. Towards a methodology for statistical disclosure control. *Statistik Tidskrift*, 15:429–444, 1977.
- [63] S. Dasgupta and A. Gupta. An elementray proof of a theorem of Johnson and Lindenstrauss. *Random Structures and Algorithms*, 22(1):60–65, 2002.

- [64] A. Dawid and V. Vovk. Prequential probability: principles and properties. *Bernoulli*, 5: 125–162, 1999.
- [65] P. Del Moral, M. Ledoux, and L. Miclo. On contraction properties of Markov kernels. *Probability Theory and Related Fields*, 126:395–420, 2003.
- [66] R. L. Dobrushin. Central limit theorem for nonstationary Markov chains. I. *Theory of Probability and Its Applications*, 1(1):65–80, 1956.
- [67] R. L. Dobrushin. Central limit theorem for nonstationary Markov chains. II. *Theory of Probability and Its Applications*, 1(4):329–383, 1956.
- [68] D. L. Donoho and J. Jin. Higher criticism for detecting sparse heterogeneous mixtures. *Annals of Statistics*, 32(3), 2004.
- [69] S. Du, S. Kakade, R. Wang, and L. Yang. Is a good representation sufficient for sample efficient reinforcement learning? In *Proceedings of the Eighth International Conference on Learning Representations*, 2020.
- [70] J. C. Duchi and R. Rogers. Lower bounds for locally private estimation via communication complexity. In *Proceedings of the Thirty Second Annual Conference on Computational Learning Theory*, 2019.
- [71] J. C. Duchi and F. Ruan. A constrained risk inequality for general losses. In *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*, 2020.
- [72] J. C. Duchi and M. J. Wainwright. Distance-based and continuum Fano inequalities with applications to statistical estimation. *arXiv:1311.2669 [cs.IT]*, 2013.
- [73] J. C. Duchi, M. I. Jordan, and H. B. McMahan. Estimation, optimization, and parallelism when data is sparse. In *Advances in Neural Information Processing Systems 26*, 2013.
- [74] J. C. Duchi, M. I. Jordan, and M. J. Wainwright. Local privacy, data processing inequalities, and minimax rates. *arXiv:1302.3203 [math.ST]*, 2013. URL <http://arxiv.org/abs/1302.3203>.
- [75] J. C. Duchi, M. I. Jordan, and M. J. Wainwright. Local privacy and statistical minimax rates. In *54th Annual Symposium on Foundations of Computer Science*, pages 429–438, 2013.
- [76] J. C. Duchi, K. Khosravi, and F. Ruan. Multiclass classification, information, divergence, and surrogate risk. *Annals of Statistics*, 46(6b):3246–3275, 2018.
- [77] R. M. Dudley. *Uniform Central Limit Theorems*. Cambridge University Press, 1999.
- [78] C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3 & 4):211–407, 2014.
- [79] C. Dwork, K. Kenthapadi, F. McSherry, I. Mironov, and M. Naor. Our data, ourselves: Privacy via distributed noise generation. In *Advances in Cryptology (EUROCRYPT 2006)*, 2006.

- [80] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *Proceedings of the Third Theory of Cryptography Conference*, pages 265–284, 2006.
- [81] C. Dwork, G. N. Rothblum, and S. P. Vadhan. Boosting and differential privacy. In *51st Annual Symposium on Foundations of Computer Science*, pages 51–60, 2010.
- [82] C. Dwork, V. Feldman, M. Hardt, T. Pitassi, O. Reingold, and A. Roth. Preserving statistical validity in adaptive data analysis. *arXiv:1411.2664v2 [cs.LG]*, 2014.
- [83] C. Dwork, V. Feldman, M. Hardt, T. Pitassi, O. Reingold, and A. Roth. Preserving statistical validity in adaptive data analysis. In *Proceedings of the Forty-Seventh Annual ACM Symposium on the Theory of Computing*, 2015.
- [84] C. Dwork, V. Feldman, M. Hardt, T. Pitassi, O. Reingold, and A. Roth. The reusable holdout: Preserving statistical validity in adaptive data analysis. *Science*, 349(6248):636–638, 2015.
- [85] G. K. Dziugaite and D. M. Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. In *Proceedings of the 33rd Conference on Uncertainty in Artificial Intelligence*, 2017.
- [86] A. V. Evfimievski, J. Gehrke, and R. Srikant. Limiting privacy breaches in privacy preserving data mining. In *Proceedings of the Twenty-Second Symposium on Principles of Database Systems*, pages 211–222, 2003.
- [87] K. Fan. Minimax theorems. *Proceedings of the National Academy of Sciences*, 39(1):42–47, 1953.
- [88] V. Feldman and T. Steinke. Calibrating noise to variance in adaptive data analysis. In *Proceedings of the Thirty First Annual Conference on Computational Learning Theory*, 2018. URL <http://arxiv.org/abs/1712.07196>.
- [89] G. Folland. *Real Analysis: Modern Techniques and their Applications*. Pure and Applied Mathematics. John Wiley & Sons, second edition, 1999.
- [90] D. Foster and R. Vohra. Asymptotic calibration. *Biometrika*, 85(2):379–390, 1998.
- [91] D. Foster, S. Kakade, J. Qian, and A. Rakhlin. The statistical complexity of interactive decision making. *arXiv:2112.13487v3 [cs.LG]*, 2021.
- [92] D. Foster, N. Golowich, and Y. Han. Tight guarantees for interactive decision making with the decision-estimation coefficient. In *Proceedings of the Thirty Sixth Annual Conference on Computational Learning Theory*, 2023.
- [93] D. J. Foster, C. Gentile, M. Mohri, and J. Zimmert. Adapting to misspecification in contextual bandits. In *Advances in Neural Information Processing Systems 33*, 2020.
- [94] D. P. Foster and S. Hart. “calibeating”: Beating forecasters at their own game. *arXiv:2209.0489 [econ.TH]*, 2022.
- [95] A. Franco, N. Malhotra, and G. Simonovits. Publication bias in the social sciences: Unlocking the file drawer. *Science*, 345(6203):1502–1505, 2014.

- [96] D. A. Freedman. On tail probabilities for martingales. *The Annals of Probability*, 3(1): 100–118, Feb. 1975.
- [97] R. Gallager. Source coding with side information and universal coding. Technical Report LIDS-P-937, MIT Laboratory for Information and Decision Systems, 1979.
- [98] D. García-García and R. C. Williamson. Divergences and risks for multiclass experiments. In *Proceedings of the Twenty Fifth Annual Conference on Computational Learning Theory*, 2012.
- [99] A. Garg, T. Ma, and H. L. Nguyen. On communication cost of distributed statistical estimation and dimensionality. In *Advances in Neural Information Processing Systems 27*, 2014.
- [100] A. Gelman and E. Loken. The garden of forking paths: Why multiple comparisons can be a problem, even when there is no “fishing expedition” or “p-hacking” and the research hypothesis was posited ahead of time. Technical report, Columbia University, 2013.
- [101] R. P. Gilbert. *Function Theoretic Methods in Partial Differential Equations*. Academic Press, 1969.
- [102] R. D. Gill and B. Y. Levit. Applications of the van Trees inequality: a Bayesian Cramér-rao bound. *Bernoulli*, 1(1–2):59–79, 1995.
- [103] T. Gneiting and A. Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007.
- [104] R. M. Gray. *Entropy and Information Theory*. Springer, 1990.
- [105] P. Grünwald. *The Minimum Description Length Principle*. MIT Press, 2007.
- [106] A. Guntuboyina. Lower bounds for the minimax risk using f -divergences, and applications. *IEEE Transactions on Information Theory*, 57(4):2386–2399, 2011.
- [107] L. Györfi and T. Nemetz. f -dissimilarity: A generalization of the affinity of several distributions. *Annals of the Institute of Statistical Mathematics*, 30:105–113, 1978.
- [108] L. Györfi, M. Kohler, A. Krzyżak, and H. Walk. *A Distribution-Free Theory of Nonparametric Regression*. Springer, 2002.
- [109] R. Z. Has’minskii. A lower bound on the risks of nonparametric estimates of densities in the uniform metric. *Theory of Probability and Applications*, 23:794–798, 1978.
- [110] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer, second edition, 2009.
- [111] J. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms I & II*. Springer, New York, 1993.
- [112] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, Mar. 1963.

- [113] N. Homer, S. Szelinger, M. Redman, D. Duggan, W. Tembe, J. Muehling, J. V. Pearson, D. A. Stephan, S. F. Nelson, and D. W. Craig. Resolving individuals contributing trace amounts of DNA to highly complex mixtures using high-density SNP genotyping microarrays. *PLoS Genetics*, 4(8):e1000167, 2008.
- [114] K. Hung and W. Fithian. Statistical methods for replicability assessment. *Annals of Applied Statistics*, 14(3):1063–1087, 2020.
- [115] I. A. Ibragimov and R. Z. Has'minskii. *Statistical Estimation: Asymptotic Theory*. Springer-Verlag, 1981.
- [116] P. Indyk. Nearest neighbors in high-dimensional spaces. In *Handbook of Discrete and Computational Geometry*. CRC Press, 2004.
- [117] P. Indyk and R. Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the Thirtieth Annual ACM Symposium on the Theory of Computing*, 1998.
- [118] J. P. Ioannidis. Why most published research findings are false. *PLoS Medicine*, 2(8), 2005. doi: 10.1371/journal.pmed.0020124.
- [119] E. T. Jaynes. On the rationale of maximum-entropy methods. *Proceedings of the IEEE*, 70(9):939–952, Sept. 1982.
- [120] T. S. Jayram. Hellinger strikes back: a note on the multi-party information complexity of AND. In *Proceedings of APPROX and RANDOM 2009*, volume 5687 of *Lecture Notes in Computer Science*, pages 562–573. Springer, 2009.
- [121] H. Jeffreys. An invariant form for the prior probability in estimation problems. *Proceedings of the Royal Society of London, Series A: Mathematical and Physical Sciences*, 186:453–461, 1946.
- [122] T. Joachims. Training linear SVMs in linear time. In *Proceedings of the ACM Conference on Knowledge Discovery and Data Mining (KDD)*, 2006.
- [123] W. Johnson and J. Lindenstrauss. Extensions of Lipschitz maps into a Hilbert space. *Contemporary Mathematics*, 26:189–206, 1984.
- [124] E. Kaufmann, N. Korda, and R. Munos. Thompson sampling: An asymptotically optimal finite-time analysis. In *Algorithmic Learning Theory*, 2012.
- [125] M. J. Kearns and L. Saul. Large deviation methods for approximate probabilistic inference. In *Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence*, pages 311–319, 1998.
- [126] A. Kolmogorov and V. Tikhomirov. ε -entropy and ε -capacity of sets in functional spaces. *Uspekhi Matematicheskikh Nauk*, 14(2):3–86, 1959.
- [127] V. Koltchinskii. *Oracle Inequalities in Empirical Risk Minimization and Sparse Recovery Problems*, volume 2033 of *Lecture Notes in Mathematics*. Springer-Verlag, 2011.
- [128] A. Kumar, P. Liang, and T. Ma. Verified uncertainty calibration. In *Advances in Neural Information Processing Systems 32*, 2019.

- [129] E. Kushilevitz and N. Nisan. *Communication Complexity*. Cambridge University Press, 1997.
- [130] T. L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6:4–22, 1985.
- [131] J. Langford and R. Caruana. (not) bounding the true error. In *Advances in Neural Information Processing Systems 14*, 2001.
- [132] T. Lattimore and C. Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- [133] T. Lattimore, C. Szepesvári, and G. Weisz. Learning with good feature representations in bandits and in RL with a generative model. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- [134] L. Le Cam. *Asymptotic Methods in Statistical Decision Theory*. Springer-Verlag, 1986.
- [135] M. Ledoux and M. Talagrand. *Probability in Banach Spaces*. Springer, 1991.
- [136] E. L. Lehmann and G. Casella. *Theory of Point Estimation, Second Edition*. Springer, 1998.
- [137] F. Liese and I. Vajda. On divergences and informations in statistics and information theory. *IEEE Transactions on Information Theory*, 52(10):4394–4412, 2006.
- [138] J. Lin. Divergence measures based on the shannon entropy. *IEEE Transactions on Information Theory*, 37(1):145–151, 1991.
- [139] J. Liu, R. van Handel, and S. Verdú. Second-order converses via reverse hypercontractivity. *Mathematical Statistics and Learning*, 2(2):103–163, 2019.
- [140] Y. Liu. Fisher consistency of multicategory support vector machines. In *Processing of 11th International Conference on Artificial Intelligence and Statistics*, pages 291–298, 2007.
- [141] D. Luenberger. *Optimization by Vector Space Methods*. Wiley, 1969.
- [142] G. Lugosi and N. Vayatis. On the Bayes-risk consistency of regularized boosting methods. *Annals of Statistics*, 32(1):30–55, 2004.
- [143] M. Madiman and A. Barron. Generalized entropy power inequalities and monotonicity properties of information. *IEEE Transactions on Information Theory*, 53(7):2317–2329, 2007.
- [144] D. A. McAllester. Some PAC-bayesian theorems. In *Proceedings of the Eleventh Annual Conference on Computational Learning Theory*, 1998.
- [145] D. A. McAllester. Simplified PAC-bayesian margin bounds. In *Proceedings of the Sixteenth Annual Conference on Computational Learning Theory*, pages 203–215, 2003.
- [146] D. A. McAllester. PAC-Bayesian stochastic model selection. *Machine Learning*, 51:5–21, 2003.
- [147] D. A. McAllester. A PAC-Bayesian tutorial with a dropout bound. *arXiv:1307.2118 [cs.LG]*, 2013.
- [148] F. McSherry and K. Talwar. Mechanism design via differential privacy. In *48th Annual Symposium on Foundations of Computer Science*, 2007.

- [149] S. Muthukrishnan. Data streams: Algorithms and applications. *Foundations and Trends in Theoretical Computer Science*, 1(2):117–236, 2005.
- [150] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.
- [151] A. Nowak-Vila, F. Bach, and A. Rudi. Consistent structured prediction with max-min margin markov networks. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- [152] D. Ostrovskii and F. Bach. Finite-sample analysis of M -estimators using self-concordance. *Electronic Journal of Statistics*, 15:326–391, 2021.
- [153] D. Petz. A survey of certain trace inequalities. *Banach Center Publications*, 30:287–298, 1994.
- [154] Y. Polyanskiy and Y. Wu. Strong data-processing inequalities for channels and Bayesian networks. In *Convexity and Concentration*, volume 161 of *The IMA Volumes in Mathematics and its Applications*. Springer, 2017.
- [155] Y. Polyanskiy and Y. Wu. *Information Theory: From Coding to Learning*. Cambridge University Press, 2024.
- [156] F. Pukelsheim. *Optimal Design of Experiments*. Classics in Applied Mathematics. SIAM, 1993.
- [157] M. Raginsky. Strong data processing inequalities and ϕ -Sobolev inequalities for discrete channels. *IEEE Transactions on Information Theory*, 62(6):3355–3389, 2016.
- [158] M. Raginsky and I. Sason. Concentration of measure inequalities in information theory, communications, and coding. *Foundations and Trends in Communications and Information Theory*, 10(1–2):1–250, 2014.
- [159] A. Rao and A. Yehudayoff. *Communication Complexity and Applications*. Cambridge University Press, 2020.
- [160] G. Raskutti, M. J. Wainwright, and B. Yu. Minimax rates of estimation for high-dimensional linear regression over ℓ_q -balls. *IEEE Transactions on Information Theory*, 57(10):6976–6994, 2011.
- [161] H. Robbins. Some aspects of the sequential design of experiments. *Bulletin American Mathematical Society*, 55:527–535, 1952.
- [162] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- [163] H. Royden. *Real Analysis*. Pearson, third edition, 1988.
- [164] D. Russo and B. Van Roy. An information-theoretic analysis of Thompson sampling. *Journal of Machine Learning Research*, page To appear, 2014.
- [165] D. Russo and B. Van Roy. Learning to optimize via information-directed sampling. In *Advances in Neural Information Processing Systems 27*, 2014.

- [166] D. Russo and B. Van Roy. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 39(4):1221–1243, 2014.
- [167] C. E. Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 1948.
- [168] A. Shapiro, D. Dentcheva, and A. Ruszczyński. *Lectures on Stochastic Programming: Modeling and Theory*. SIAM and Mathematical Programming Society, 2009.
- [169] G. R. Shorack and J. A. Wellner. *Empirical processes with applications to statistics*, volume 59. Siam, 2009.
- [170] M. Sion. On general minimax theorems. *Pacific Journal of Mathematics*, 8(1):171–176, 1958.
- [171] A. Slavkovic and F. Yu. Genomics and privacy. *Chance*, 28(2):37–39, 2015.
- [172] C. Stein. Efficient nonparametric testing and estimation. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, pages 187–195, 1956.
- [173] I. Steinwart. How to compare different loss functions and their risks. *Constructive Approximation*, 26:225–287, 2007.
- [174] C. Sudlow, J. Gallacher, N. Allen, V. Beral, P. Burton, J. Danesh, P. Downey, P. Elliott, J. Green, M. Landray, B. Liu, P. Matthews, G. Ong, J. Pell, A. Silman, A. Young, T. Sprosen, T. Peakman, and R. Collins. UK Biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Medicine*, 12(3):e1001779, 2015.
- [175] R. Sutton and A. Barto. *Reinforcement Learning: An Introduction (Second Edition)*. MIT Press, 2018.
- [176] T. Tao. *An Epsilon of Room, I: Real Analysis (pages from year three of a mathematical blog)*, volume 117 of *Graduate Studies in Mathematics*. American Mathematical Society, 2010.
- [177] B. Taskar. *Learning Structured Prediction Models: A Large Margin Approach*. PhD thesis, Stanford University, 2005.
- [178] B. Taskar, D. Klein, M. Collins, D. Koller, and C. Manning. Max-margin parsing. In *Empirical Methods in Natural Language Processing*, 2004.
- [179] W. R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- [180] R. Tibshirani, J. Taylor, R. Lockhart, and R. Tibshirani. Exact post-selection inference for sequential regression procedures. *Journal of the American Statistical Association*, 111(514):600–620, 2016.
- [181] J. A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12:389–434, 2012.
- [182] A. B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, 2009.
- [183] J. W. Tukey. *Exploratory Data Analysis*. Pearson, 1997.

- [184] A. W. van der Vaart. Superefficiency. In D. Pollard, E. Torgersen, and G. Yang, editors, *Festschrift for Lucien Le Cam*, chapter 27. Springer, 1997.
- [185] A. W. van der Vaart. *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1998.
- [186] A. W. van der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer, New York, 1996.
- [187] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. In *Compressed Sensing: Theory and Applications*, chapter 5, pages 210–268. Cambridge University Press, 2012.
- [188] C. Villani. *Optimal Transport: Old and New*. Springer, 2009.
- [189] A. Wald. Contributions to the theory of statistical estimation and testing hypotheses. *Annals of Mathematical Statistics*, 10(4):299–326, 1939.
- [190] S. Warner. Randomized response: a survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60(309):63–69, 1965.
- [191] Y. Wu and P. Yang. Minimax rates of entropy estimation on large alphabets via best polynomial approximation. *IEEE Transactions on Information Theory*, 62(6):3702–3720, 2016.
- [192] Y. Yang and A. Barron. Information-theoretic determination of minimax rates of convergence. *Annals of Statistics*, 27(5):1564–1599, 1999.
- [193] A. C.-C. Yao. Some complexity questions related to distributive computing (preliminary report). In *Proceedings of the Eleventh Annual ACM Symposium on Theory of Computing*, pages 209–213. ACM, 1979.
- [194] B. Yu. Assouad, Fano, and Le Cam. In *Festschrift for Lucien Le Cam*, pages 423–435. Springer-Verlag, 1997.
- [195] T. Zhang. Statistical analysis of some multi-category large margin classification methods. *Journal of Machine Learning Research*, 5:1225–1251, 2004.
- [196] T. Zhang. Statistical behavior and consistency of classification methods based on convex risk minimization. *Annals of Statistics*, 32:56–85, 2004.
- [197] T. Zhang and B. Yu. Boosting with early stopping: Convergence and consistency. *The Annals of Statistics*, 33:1538–1579, 2005.
- [198] Y. Zhang, J. C. Duchi, M. I. Jordan, and M. J. Wainwright. Information-theoretic lower bounds for distributed estimation with communication constraints. In *Advances in Neural Information Processing Systems 26*, 2013.
- [199] N. Zhivotovskiy. Dimension-free bounds for sums of independent matrices and simple tensors via the variational principle. *Electronic Journal of Probability*, 29:1–28, 2024.