

The Elements of Differentiable Programming

Mathieu Blondel

Google DeepMind
mblondel@google.com

Vincent Roulet

Google DeepMind
vroulet@google.com

Contents

1	Introduction	6
1.1	What is differentiable programming?	6
1.2	Book goals and scope	10
1.3	Intended audience	11
1.4	How to read this book?	11
1.5	Related work	11
I	Fundamentals	13
2	Differentiation	14
2.1	Univariate functions	14
2.1.1	Derivatives	14
2.1.2	Calculus rules	18
2.1.3	Leibniz's notation	20
2.2	Multivariate functions	21
2.2.1	Directional derivatives	21
2.2.2	Gradients	22
2.2.3	Jacobians	26
2.3	Linear maps	31
2.3.1	The need for linear maps	31
2.3.2	Euclidean spaces	32

2.3.3	Linear maps and their adjoints	33
2.3.4	Jacobian-vector products	34
2.3.5	Vector-Jacobian products	36
2.3.6	Chain rule using linear maps	37
2.3.7	Functions of multiple inputs (fan-in)	38
2.3.8	Functions with multiple outputs (fan-out)	41
2.3.9	Extensions to non-Euclidean linear spaces	41
2.4	Second-order differentiation	43
2.4.1	Second derivatives	43
2.4.2	Second directional derivatives	44
2.4.3	Hessians	44
2.4.4	Hessian-vector products	46
2.4.5	Second-order Jacobians	46
2.5	Higher-order differentiation	47
2.5.1	Higher-order derivatives	47
2.5.2	Higher-order directional derivatives	47
2.5.3	Higher-order Jacobians	48
2.5.4	Taylor expansions	49
2.6	Differential geometry	50
2.6.1	Differentiability on manifolds	50
2.6.2	Tangent spaces and pushforward operators	51
2.6.3	Cotangent spaces and pullback operators	52
2.7	Generalized derivatives	55
2.7.1	Rademacher's theorem	56
2.7.2	Clarke derivatives	56
2.8	Summary	58
3	Probabilistic learning	61
3.1	Probability distributions	61
3.1.1	Discrete probability distributions	61
3.1.2	Continuous probability distributions	62
3.2	Maximum likelihood estimation	63
3.2.1	Negative log-likelihood	63
3.2.2	Consistency w.r.t. the Kullback-Leibler divergence	63
3.3	Probabilistic supervised learning	64
3.3.1	Conditional probability distributions	64

3.3.2	Inference	64
3.3.3	Binary classification	65
3.3.4	Multiclass classification	67
3.3.5	Regression	69
3.3.6	Multivariate regression	70
3.3.7	Integer regression	71
3.3.8	Loss functions	72
3.4	Exponential family distributions	75
3.4.1	Definition	75
3.4.2	The log-partition function	76
3.4.3	Maximum entropy principle	78
3.4.4	Maximum likelihood estimation	79
3.4.5	Probabilistic learning with exponential families	80
3.5	Summary	81

II Differentiable programs 83

4 Parameterized programs 84

4.1	Representing computer programs	84
4.1.1	Computation chains	84
4.1.2	Directed acyclic graphs	85
4.1.3	Computer programs as DAGs	87
4.1.4	Arithmetic circuits	89
4.2	Feedforward networks	90
4.3	Multilayer perceptrons	90
4.3.1	Combining affine layers and activations	90
4.3.2	Link with generalized linear models	91
4.4	Activation functions	92
4.4.1	ReLU and softplus	92
4.4.2	Max pooling and log-sum-exp	92
4.4.3	Sigmoids: binary step and logistic functions	94
4.4.4	Probability mappings: argmax and softargmax	96
4.5	Normalization layers	97
4.5.1	Batch normalization	98
4.5.2	Layer normalization	99

4.6	Residual neural networks	100
4.7	Recurrent neural networks	101
4.7.1	Vector to sequence	101
4.7.2	Sequence to vector	103
4.7.3	Sequence to sequence (aligned)	103
4.7.4	Sequence to sequence (unaligned)	104
4.8	Transformers	105
4.8.1	Attention	105
4.8.2	Self-attention	106
4.8.3	Multi-head attention	107
4.8.4	Transformer layer	108
4.8.5	Transformer block	109
4.8.6	Token encoding	110
4.8.7	Positional encoding	111
4.8.8	Decoder-only architectures	116
4.8.9	Encoder-only architectures	119
4.8.10	Encoder-decoder architectures	120
4.9	Summary	122
5	Control flows	123
5.1	Comparison operators	123
5.2	Soft inequality operators	125
5.2.1	Heuristic definition	125
5.2.2	Stochastic process perspective	126
5.3	Soft equality operators	129
5.3.1	Heuristic definition	129
5.3.2	Stochastic process perspective	130
5.3.3	Gaussian process perspective	133
5.4	Logical operators	134
5.5	Continuous extensions of logical operators	135
5.5.1	Probabilistic continuous extension	135
5.5.2	Triangular norms and co-norms	137
5.6	If-else statements	138
5.6.1	Differentiating through branch variables	139
5.6.2	Differentiating through predicate variables	140
5.6.3	Continuous relaxations	141

5.7	Else-if statements	144
5.7.1	Encoding K branches	144
5.7.2	Conditionals	145
5.7.3	Differentiating through branch variables	146
5.7.4	Differentiating through predicate variables	147
5.7.5	Continuous relaxations	148
5.8	For loops	149
5.9	Scan functions	151
5.10	While loops	152
5.10.1	While loops as cyclic graphs	152
5.10.2	Unrolled while loops	154
5.10.3	Markov chain perspective	157
5.11	Summary	159
6	Data structures	161
6.1	Lists	161
6.1.1	Basic operations	162
6.1.2	Operations on variable-length lists	163
6.1.3	Continuous relaxations using soft indexing	165
6.2	Dictionaries	168
6.2.1	Basic operations	168
6.2.2	Continuous relaxation using kernel regression	170
6.2.3	Discrete probability distribution perspective	171
6.2.4	Link with attention in Transformers	172
6.3	Summary	173
III	Differentiating through programs	175
7	Finite differences	176
7.1	Forward differences	176
7.2	Backward differences	177
7.3	Central differences	178
7.4	Higher-accuracy finite differences	179
7.5	Higher-order finite differences	180
7.6	Complex-step derivatives	181

7.7	Complexity	182
7.8	Summary	182
8	Automatic differentiation	184
8.1	Computation chains	184
8.1.1	Forward-mode	185
8.1.2	Reverse-mode	188
8.1.3	Complexity of computing entire Jacobians	192
8.2	Feedforward networks	194
8.2.1	Computing the adjoint	194
8.2.2	Computing the gradient	195
8.3	Computation graphs	197
8.3.1	Forward mode	197
8.3.2	Reverse mode	200
8.3.3	Complexity, the Baur-Strassen theorem	204
8.4	Implementation	204
8.4.1	Primitive functions	204
8.4.2	Closure under function composition	205
8.4.3	Examples of JVPs and VJPs	205
8.4.4	Automatic linear transposition	206
8.5	Checkpointing	207
8.5.1	Recursive halving	209
8.5.2	Dynamic programming	211
8.5.3	Online checkpointing	213
8.6	Reversible layers	213
8.6.1	General case	213
8.6.2	Case of orthonormal JVPs	214
8.7	Randomized forward-mode gradient estimator	215
8.8	Summary	216
9	Second-order automatic differentiation	217
9.1	Hessian-vector products	217
9.1.1	Four possible methods	217
9.1.2	Complexity	218
9.2	Gauss-Newton matrix	222
9.2.1	An approximation of the Hessian	222

9.2.2	Gauss-Newton chain rule	223
9.2.3	Gauss-Newton vector product	223
9.2.4	Gauss-Newton matrix factorization	224
9.2.5	Stochastic setting	225
9.3	Fisher information matrix	225
9.3.1	Definition using the score function	225
9.3.2	Link with the Hessian	226
9.3.3	Equivalence with the Gauss-Newton matrix	226
9.4	Inverse-Hessian vector product	228
9.4.1	Definition as a linear map	228
9.4.2	Implementation with matrix-free linear solvers	228
9.4.3	Complexity	229
9.5	Second-order backpropagation	230
9.5.1	Second-order Jacobian chain rule	230
9.5.2	Computation chains	232
9.5.3	Fan-in and fan-out	233
9.6	Block diagonal approximations	234
9.6.1	Feedforward networks	234
9.6.2	Computation graphs	236
9.7	Diagonal approximations	236
9.7.1	Computation chains	237
9.7.2	Computation graphs	238
9.8	Randomized estimators	239
9.8.1	Girard-Hutchinson estimator	239
9.8.2	Bartlett estimator for the factorization	240
9.8.3	Bartlett estimator for the diagonal	241
9.9	Summary	242
10	Inference in graphical models as differentiation	243
10.1	Chain rule of probability	243
10.2	Conditional independence	244
10.3	Inference problems	245
10.3.1	Joint probability distributions	245
10.3.2	Likelihood	245
10.3.3	Maximum a-posteriori inference	246
10.3.4	Marginal inference	246

10.3.5	Expectation, convex hull, marginal polytope	247
10.3.6	Complexity of brute force	248
10.4	Markov chains	248
10.4.1	The Markov property	249
10.4.2	Time-homogeneous Markov chains	251
10.4.3	Higher-order Markov chains	252
10.5	Bayesian networks	252
10.5.1	Expressing variable dependencies using DAGs	252
10.5.2	Parameterizing Bayesian networks	253
10.5.3	Ancestral sampling	254
10.6	Markov random fields	254
10.6.1	Expressing factors using undirected graphs	254
10.6.2	MRFs as exponential family distributions	255
10.6.3	Conditional random fields	257
10.6.4	Sampling	257
10.7	Inference on chains	257
10.7.1	The forward-backward algorithm	258
10.7.2	The Viterbi algorithm	259
10.8	Inference on trees	261
10.9	Inference as differentiation	262
10.9.1	Inference as gradient of the log-partition	262
10.9.2	Semirings and softmax operators	263
10.9.3	Inference as backpropagation	265
10.10	Summary	267
11	Differentiating through optimization	269
11.1	Implicit functions	269
11.1.1	Optimization problems	270
11.1.2	Nonlinear equations	270
11.1.3	Application to bilevel optimization	270
11.2	Envelope theorems	271
11.2.1	Danskin's theorem	272
11.2.2	Rockafellar's theorem	273
11.3	Implicit function theorem	274
11.3.1	Univariate functions	274
11.3.2	Multivariate functions	276

11.3.3	JVP and VJP of implicit functions	278
11.3.4	Proof of the implicit function theorem	279
11.4	Adjoint state method	280
11.4.1	Differentiating nonlinear equations	280
11.4.2	Relation with envelope theorems	281
11.4.3	Proof using the method of Lagrange multipliers	281
11.4.4	Proof using the implicit function theorem	282
11.4.5	Reverse mode as adjoint method with backsubstitution	282
11.5	Inverse function theorem	285
11.5.1	Differentiating inverse functions	285
11.5.2	Link with the implicit function theorem	286
11.5.3	Proof of inverse function theorem	286
11.6	Summary	287
12	Differentiating through integration	289
12.1	Differentiation under the integral sign	289
12.2	Differentiating through expectations	290
12.2.1	Parameter-independent distributions	290
12.2.2	Parameter-dependent distributions	291
12.2.3	Application to expected loss functions	293
12.2.4	Application to experimental design	294
12.3	Score function estimators, REINFORCE	295
12.3.1	Scalar-valued functions	295
12.3.2	Variance reduction	298
12.3.3	Vector-valued functions	299
12.3.4	Second derivatives	300
12.4	Path gradient estimators, reparametrization trick	301
12.4.1	Location-scale transforms	301
12.4.2	Differentiable transforms	303
12.4.3	Inverse transforms	304
12.4.4	Pushforward operators	306
12.4.5	Change-of-variables theorem	308
12.5	Stochastic programs	309
12.5.1	Stochastic computation graphs	309
12.5.2	Examples	312
12.5.3	Unbiased gradient estimators	314

12.5.4	Local vs. global expectations	316
12.6	Differential equations	317
12.6.1	Parameterized differential equations	317
12.6.2	Continuous adjoint method	320
12.6.3	Gradients via the continuous adjoint method	321
12.6.4	Gradients via reverse-mode on discretization	324
12.6.5	Reversible discretization schemes	325
12.6.6	Proof of the continuous adjoint method	328
12.7	Summary	329

IV Smoothing programs 332

13 Smoothing by optimization 333

13.1	Primal approach	333
13.1.1	Infimal convolution	334
13.1.2	Moreau envelope	335
13.1.3	Vector-valued functions	339
13.2	Legendre–Fenchel transforms, convex conjugates	341
13.2.1	Definition	341
13.2.2	Closed-form examples	342
13.2.3	Properties	344
13.2.4	Conjugate calculus	346
13.2.5	Fast Legendre transform	346
13.3	Dual approach	347
13.3.1	Duality between strong convexity and smoothness	347
13.3.2	Smoothing by dual regularization	348
13.3.3	Equivalence between primal and dual regularizations	350
13.3.4	Regularization scaling	351
13.3.5	Generalized entropies	352
13.4	Smoothed ReLU functions	356
13.5	Smoothed max operators	358
13.5.1	Definition and properties	358
13.5.2	Reduction to root finding	359
13.5.3	The softmax	360
13.5.4	The sparsemax	361

13.5.5 Recovering smoothed ReLU functions	364
13.6 Relaxed step functions (sigmoids)	364
13.7 Relaxed argmax operators	365
13.8 Summary	369
14 Smoothing by integration	371
14.1 Convolution	371
14.1.1 Convolution operators	371
14.1.2 Convolution with a kernel	372
14.1.3 Discrete convolution	373
14.1.4 Differentiation	375
14.1.5 Multidimensional convolution	375
14.1.6 Link between convolution and infimal convolution	375
14.1.7 The soft infimal convolution	376
14.1.8 The soft Moreau envelope	377
14.2 Fourier and Laplace transforms	378
14.2.1 Convolution theorem	378
14.2.2 Link between Fourier and Legendre transforms	378
14.2.3 The soft Legendre-Fenchel transform	379
14.3 Examples	382
14.3.1 Smoothed step function	382
14.3.2 Smoothed ReLU function	383
14.4 Perturbation of blackbox functions	385
14.4.1 Expectation in a location-scale family	385
14.4.2 Gradient estimation by reparametrization	386
14.4.3 Gradient estimation by SFE, Stein's lemma	387
14.4.4 Link between reparametrization and SFE	388
14.4.5 Variance reduction and evolution strategies	389
14.4.6 Zero-temperature limit	390
14.5 Gumbel tricks	391
14.5.1 The Gumbel distribution	391
14.5.2 Perturbed comparison	393
14.5.3 Perturbed argmax	394
14.5.4 Perturbed max	395
14.5.5 Gumbel trick for sampling	396
14.5.6 Perturb-and-MAP	396

14.5.7	Gumbel-softargmax	398
14.6	Summary	400
V	Optimizing differentiable programs	402
15	Optimization basics	403
15.1	Objective functions	403
15.2	Oracles	404
15.3	Variational perspective of optimization algorithms	405
15.4	Classes of functions	405
15.4.1	Lipschitz functions	405
15.4.2	Smooth functions	406
15.4.3	Convex functions	408
15.4.4	Strongly-convex functions	410
15.4.5	Nonconvex functions	411
15.5	Performance guarantees	413
15.6	Summary	416
16	First-order optimization	417
16.1	Gradient descent	417
16.1.1	Variational perspective	417
16.1.2	Convergence for smooth functions	418
16.1.3	Momentum and accelerated variants	420
16.2	Stochastic gradient descent	421
16.2.1	Stochastic gradients	422
16.2.2	Vanilla SGD	423
16.2.3	Momentum variants	424
16.2.4	Adaptive variants	425
16.3	Projected gradient descent	425
16.3.1	Variational perspective	426
16.3.2	Optimality conditions	427
16.3.3	Commonly-used projections	427
16.4	Proximal gradient method	428
16.4.1	Variational perspective	429
16.4.2	Optimality conditions	429

16.4.3	Commonly-used proximal operators	430
16.5	Summary	430
17	Second-order optimization	432
17.1	Newton's method	432
17.1.1	Variational perspective	432
17.1.2	Regularized Newton method	433
17.1.3	Approximate direction	434
17.1.4	Convergence guarantees	434
17.1.5	Linesearch	434
17.1.6	Geometric interpretation	435
17.1.7	Stochastic Newton's method	436
17.2	Gauss-Newton method	437
17.2.1	With exact outer function	438
17.2.2	With approximate outer function	439
17.2.3	Linesearch	440
17.2.4	Stochastic Gauss-Newton	440
17.3	Natural gradient descent	441
17.3.1	Variational perspective	441
17.3.2	Stochastic natural gradient descent	442
17.4	Quasi-Newton methods	443
17.4.1	BFGS	443
17.4.2	Limited-memory BFGS	444
17.5	Approximate Hessian diagonal inverse preconditioners	444
17.6	Summary	444
18	Duality	446
18.1	Dual norms	446
18.2	Fenchel duality	447
18.3	Bregman divergences	450
18.4	Fenchel-Young loss functions	453
18.5	Summary	454
	References	455

The Elements of Differentiable Programming

Mathieu Blondel¹ and Vincent Roulet¹

¹*Google DeepMind*

ABSTRACT

Artificial intelligence has recently experienced remarkable advances, fueled by large models, vast datasets, accelerated hardware, and, last but not least, the transformative power of differentiable programming. This new programming paradigm enables end-to-end differentiation of complex computer programs (including those with control flows and data structures), making gradient-based optimization of program parameters possible.

As an emerging paradigm, differentiable programming builds upon several areas of computer science and applied mathematics, including automatic differentiation, graphical models, optimization and statistics. This book presents a comprehensive review of the fundamental concepts useful for differentiable programming. We adopt two main perspectives, that of optimization and that of probability, with clear analogies between the two.

Differentiable programming is not merely the differentiation of programs, but also the thoughtful design of programs intended for differentiation. By making programs differentiable, we inherently introduce probability distributions over their execution, providing a means to quantify the uncertainty associated with program outputs.

Acknowledgements

We thank the following people for sending us feedback, suggestions and typos: Fabian Pedregosa, Kevin Murphy, Niklas Schmitz, Nidham Gazagnadou, Bruno De Backer, David López, Guillaume Gautier, Sam Duffield, Logan Bruns, Wojciech Stokowiec, Alex Towell, John Reid, Sadish Dhakal, Fabian Schaipp, Mahmoud Asem, Simone Scardapane, (add your name here!).

Source code

We provide some Python source code to accompany the book on [github](#).

Notation

Table 1: Naming conventions

Notation	Description
$\mathcal{X} \subseteq \mathbb{R}^D$	Input space (e.g., features)
$\mathcal{Y} \subseteq \mathbb{R}^M$	Output space (e.g., classes)
$\mathcal{S}_k \subseteq \mathbb{R}^{D_k}$	Output space on layer or state k
$\mathcal{W} \subseteq \mathbb{R}^P$	Weight space
$\Lambda \subseteq \mathbb{R}^Q$	Hyperparameter space
$\Theta \subseteq \mathbb{R}^R$	Distribution parameter space, logit space
N	Number of training samples
T	Number of optimization iterations
$\boldsymbol{x} \in \mathcal{X}$	Input vector
$\boldsymbol{y} \in \mathcal{Y}$	Target vector
$\boldsymbol{s}_k \in \mathcal{S}_k$	State vector k
$\boldsymbol{w} \in \mathcal{W}$	Network (model) weights
$\boldsymbol{\lambda} \in \Lambda$	Hyperparameters
$\boldsymbol{\theta} \in \Theta$	Distribution parameters, logits
$\pi \in [0, 1]$	Probability value
$\boldsymbol{\pi} \in \Delta^M$	Probability vector

Table 2: Naming conventions (continued)

Notation	Description
f	Network function
$f(\cdot; \mathbf{x})$	Network function with \mathbf{x} fixed
L	Objective function
ℓ	Loss function
κ	Kernel function
ϕ	Output embedding, sufficient statistic
step	Heaviside step function
logistic $_{\sigma}$	Logistic function with temperature σ
logistic	Shorthand for logistic $_1$
$p_{\boldsymbol{\theta}}$	Model distribution with parameters $\boldsymbol{\theta}$
ρ	Data distribution over $\mathcal{X} \times \mathcal{Y}$
$\rho_{\mathcal{X}}$	Data distribution over \mathcal{X}
$\boldsymbol{\mu}, \sigma^2$	Mean and variance
Z	Random noise variable

1

Introduction

1.1 What is differentiable programming?

A computer program is a sequence of elementary instructions for performing a task. In traditional computer programming, the program is typically manually written by a programmer. However, for certain tasks, particularly those involving intricate patterns and complex decision-making, such as image recognition or text generation, manually writing a program is extremely challenging, if not impossible.

In contrast, modern neural networks offer a different approach. They are constructed by combining parameterized functional blocks and are trained directly from data using gradient-based optimization. This end-to-end training process, where the network learns both feature extraction and task execution simultaneously, allows neural networks to tackle complex tasks that were previously considered insurmountable for traditional, hand-coded programs. This new programming paradigm has been referred to as “differentiable programming” or “software 2.0”. We give an informal definition below.

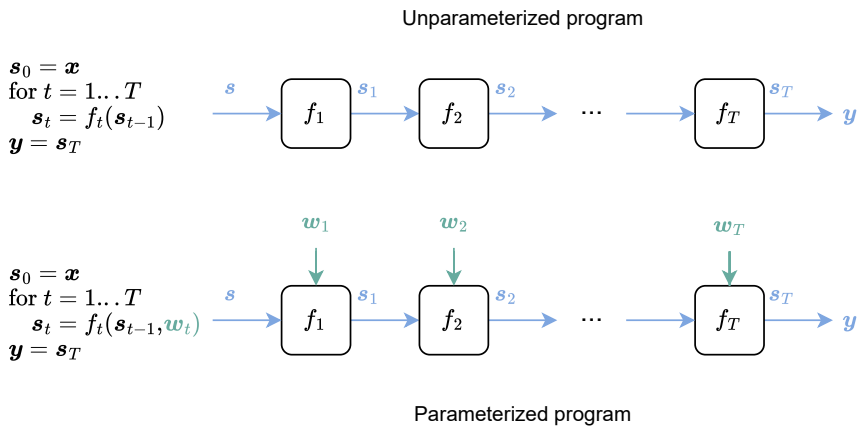


Figure 1.1: Neural networks can be seen as parameterized programs.

Definition 1.1 (Differentiable programming). Differentiable programming is a programming paradigm in which complex computer programs (including those with control flows and data structures) can be differentiated end-to-end automatically, enabling gradient-based optimization of parameters in the program.

Neural networks as parameterized programs

In differentiable programming, as in regular computer programming, a program is defined as the composition of elementary operations, forming a **computation graph**. The key difference is that, as illustrated in Fig. 1.1, the program (such as a neural network) contains **parameters** that can be adjusted from data and can be differentiated end-to-end, using **automatic differentiation** (autodiff). Typically, it is assumed that the program defines a **mathematically valid function** (a.k.a. pure function): the function should return identical values for identical arguments and should not have any side effects. Moreover, the function should have **well-defined derivatives**, ensuring that it can be used in a gradient-based optimization algorithm. Therefore, differentiable programming is not only the art of differentiating through programs but also of **designing** meaningful differentiable programs.

Why do we need derivatives?

Machine learning typically boils down to optimizing a certain objective function, which is the composition of a loss function and a model (network) function. Derivative-free optimization is called **zero-order optimization**. It only assumes that we can evaluate the objective function that we wish to optimize. Unfortunately, it is known to suffer from the **curse of dimensionality**, i.e., it only scales to small dimensional problems, such as less than 10 dimensions. Derivative-based optimization, on the other hand, is much more efficient and can scale to millions or billions of parameters. Algorithms that use first and second derivatives are known as **first-order** and **second-order** algorithms, respectively.

Why is autodiff so useful?

Before the autodiff revolution, researchers and practitioners needed to manually implement the gradient of the functions they wished to optimize. Manually deriving gradients can become very tedious for complicated functions. Moreover, every time the function is changed (for example, for trying out a new idea), the gradient needs to be re-derived. Autodiff is a game changer because it allows users to focus on quickly and creatively experimenting with functions for their tasks. An example of JAX code (Bradbury *et al.*, 2018) is given in Fig. 1.2.

Differentiable programming is not just deep learning

While there is clearly overlap between deep learning and differentiable programming, their focus is different. Deep learning studies artificial neural networks composed of multiple layers, able to learn **intermediate representations** of the data. Neural network architectures have been proposed with various **inductive biases**. For example, convolutional neural networks are designed for images and transformers are designed for sequences. On the other hand, differentiable programming studies the techniques for designing complex programs and differentiating through them. It is useful beyond deep learning: for instance in reinforcement learning, probabilistic programming and scientific computing in general.

```
import jax.numpy as jnp
from jax import grad, jit

def predict(params, inputs):
    for W, b in params:
        outputs = jnp.dot(inputs, W) + b
        inputs = jnp.tanh(outputs)
    return outputs

def loss_fn(params, inputs, targets):
    outputs = predict(params, inputs)
    return jnp.sum((outputs - targets) ** 2)

grad_fun = jit(grad(loss_fn))
```

Figure 1.2: Thanks to automatic differentiation (autodiff), the user can focus on expressing the forward computation (model), enabling fast experimentation and alleviating the need for error-prone manual gradient derivation.

Differentiable programming is not just autodiff

While autodiff is a key ingredient of differentiable programming, this is not the only one. Differentiable programming is also concerned with the design of principled differentiable operations. In fact, much research on differentiable programming has been devoted to make classical computer programming operations compatible with autodiff. As we shall see, many differentiable relaxations can be interpreted in a probabilistic framework. A core theme of this book is the interplay between optimization, probability and differentiation. Differentiation is useful for optimization and conversely, optimization can be used to design differentiable operators.

Our vision for differentiable programming

Computer programming offers powerful tools like control flows, data structures, and standard libraries, enabling users to construct complex programs for solving intricate problems. Our long-term vision is to achieve parity between traditional and differentiable programming, empowering programmers to seamlessly express differentiable programs (such as neural networks) using the full suite of tools they are accus-

tomed to. However, as discussed earlier, differentiable programming is not simply a matter of applying automatic differentiation to existing code. Programs must be designed with differentiability in mind. This usually comes to inducing a probability distribution over the program or its components. While significant work remains to fully realize this ambitious goal, we hope this book offers a solid foundation.

Where does “differentiable programming” come from?

While neural networks and autodiff have existed for several decades, the term “differentiable programming” is more recent. Olah (2015) discussed analogies between neural network architectures and higher-order functions in functional programming, and referred to it as a “new kind of programming”. Dalrymple (2016) wrote an essay titled “differentiable programming”, recognizing a paradigm where programs learn details through differentiation, with the expressiveness of functional programming. Plotkin (2018) gave a keynote talk titled “Some principles of differential programming languages” at POPL 2018. The “differentiable programming” and “software 2.0” terms were popularized among others by LeCun (2018) and Karpathy (2017). From this perspective, autodiff frameworks can be seen, not merely as libraries, but as domain-specific languages (DSLs) embedded into an existing programming language, such as Python. See also Imai (2019) for a review.

1.2 Book goals and scope

The present book aims to provide an comprehensive introduction to differentiable programming with an emphasis on **core** mathematical tools.

- In Part I, we review **fundamentals**: differentiation and probabilistic learning.
- In Part II, we review **differentiable programs**. This includes neural networks, sequence networks and control flows.
- In Part III, we review how to **differentiate through programs**. This includes automatic differentiation, but also differentiating

through optimization and integration (in particular, expectations).

- In Part [IV](#), we review **smoothing programs**. We focus on two main techniques: infimal convolution, which comes from the world of optimization and convolution, which comes from the world of integration. We also strive to spell out the connections between them.
- In Part [V](#), we review **optimizing programs**: basic optimization concepts, first-order algorithms, second-order algorithms and duality.

Our goal is to present the fundamental techniques useful for differentiable programming, **not** to survey how these techniques have been used in various applications.

1.3 Intended audience

This book is intended to be a graduate-level introduction to differentiable programming. Our pedagogical choices are made with the machine learning community in mind. Some familiarity with calculus, linear algebra, probability theory and machine learning is beneficial.

1.4 How to read this book?

This book does not need to be read linearly chapter by chapter. When needed, we indicate at the beginning of a chapter what chapters are recommended to be read as a prerequisite.

1.5 Related work

Differentiable programming builds upon a variety of connected topics. We review in this section relevant textbooks, tutorials and software.

Standard textbooks on backpropagation and automatic differentiation (autodiff) are that of Werbos ([1994](#)) and Griewank and Walther ([2008](#)). A tutorial with a focus on machine learning is provided by Baydin *et al.* ([2018](#)). Automatic differentiation is also reviewed as part

of more general textbooks, such as those of Deisenroth *et al.* (2020), Murphy (2022) (from a linear algebra perspective) and Murphy (2023) (from a functional perspective; autodiff section authored by Roy Frostig). The present book was also influenced by Peyré (2020)’s textbook on data science. The history of reverse-mode autodiff is reviewed by Griewank (2012).

A tutorial on different perspectives of backpropagation is “There and Back Again: A Tale of Slopes and Expectations” (link), by Deisenroth and Ong. A tutorial on implicit differentiation is “Deep Implicit Layers - Neural ODEs, Deep Equilibrium Models, and Beyond” (link), by Kolter, Duvenaud, and Johnson.

The standard reference on inference in graphical models and its connection with exponential families is that of Wainwright and Jordan (2008). Differential programming is also related to probabilistic programming; see, e.g., Meent *et al.* (2018).

A review of smoothing from the infimal convolution perspective is provided by Beck and Teboulle (2012). A standard textbook on convex optimization is that of Nesterov (2018). A textbook on first-order optimization methods is that of Beck (2017).

Autodiff implementations that accelerated the autodiff revolution in machine learning are Theano (Bergstra *et al.*, 2010) and Autograd (Maclaurin *et al.*, 2015). Major modern implementations of autodiff include Tensorflow (Abadi *et al.*, 2016), JAX (Bradbury *et al.*, 2018), and PyTorch (Paszke *et al.*, 2019). We in particular acknowledge the JAX team for influencing our view of autodiff.

Part I

Fundamentals

2

Differentiation

In this chapter, we review key differentiation concepts. In particular, we emphasize on the fundamental role played by linear maps.

2.1 Univariate functions

2.1.1 Derivatives

To study functions, we need to capture their infinitesimal variations around points as defined by the notion of **limit**.

Definition 2.1 (Limit). We say that $c \in \mathbb{R}$ is the **limit** of $f: \mathbb{R} \rightarrow \mathbb{R}$ as $v \in \mathbb{R}$ approaches $w \in \mathbb{R}$, denoted

$$\lim_{v \rightarrow w} f(v) = c,$$

if, for any $\varepsilon > 0$, there exists $R > 0$ such that for any $v \in \mathbb{R}$ satisfying $0 < |v - w| \leq R$, we have $|f(v) - c| \leq \varepsilon$.

We can also write $f(v) \rightarrow c$ as $v \rightarrow w$. Limits are preserved under additions and multiplications. Namely, if $\lim_{v \rightarrow w} f(v) = c$ and $\lim_{v \rightarrow w} g(v) = d$, then denoting $(af + bg)(w) := af(w) + bg(w)$ for any $a, b \in \mathbb{R}$ and $(fg)(w) := f(w)g(w)$, we have by definition of the

limit, $\lim_{v \rightarrow w}(af + bg)(v) = ac + bd$ and $\lim_{v \rightarrow w}(fg)(v) = cd$. The preservation of the limit under addition and multiplication by a scalar is generally referred to as the linearity of the limit, a property that many definitions in the sequel inherit.

With the notion of limit, we can already delineate a class of “well-behaved” functions: functions whose limits at any point equals the value of the function at that point. Functions satisfying this property are called **continuous**.

Definition 2.2 (Continuous function). A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is continuous at a point $w \in \mathbb{R}$ if

$$\lim_{v \rightarrow w} f(v) = f(w).$$

A function f is said to be continuous if it is continuous at all points in its domain.

Although the notion of continuity appears to be a benign assumption, several functions commonly-used in machine learning, such as the Heaviside step function (displayed in the left panel of Fig. 2.2), are not continuous and require special treatment.

Remark 2.1 (Little o notation). In the following, we will make use of Landau’s little o notation. We write

$$g(v) = o(f(v)) \text{ as } v \rightarrow w$$

if

$$\lim_{v \rightarrow w} \frac{|g(v)|}{|f(v)|} = 0.$$

That is, the function f dominates g in the limit $v \rightarrow w$. For example, f is continuous at w if and only if

$$f(w + \delta) = f(w) + o(1) \text{ as } \delta \rightarrow 0.$$

We now explain derivatives. Consider a function $f : \mathbb{R} \rightarrow \mathbb{R}$. As illustrated in Fig. 2.1, its value on an interval $[w_0, w_0 + \delta]$ can be approximated by the secant between its values $f(w_0)$ and $f(w_0 + \delta)$, a linear function with slope $(f(w_0 + \delta) - f(w_0))/\delta$. In the limit of an

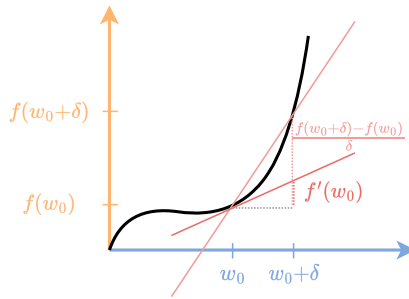


Figure 2.1: A function f can be locally approximated around a point w_0 by a secant, a linear function $w \mapsto aw + b$ with slope a and intercept b , crossing f at w_0 with value $u_0 = f(w_0)$ and crossing at $w_0 + \delta$ with value $u_\delta = f(w_0 + \delta)$. Using $u_0 = aw_0 + b$ and $u_\delta = a(w_0 + \delta) + b$, we find that its slope is $a = (f(w_0 + \delta) - f(w_0))/\delta$ and the intercept is $b = f(w_0) - aw_0$. The derivative $f'(w)$ of a function f at a point w_0 is then defined as the limit of the slope a when $\delta \rightarrow 0$. It is the slope of the tangent of f at w_0 . The value $f(w)$ of the function at w can then be locally approximated around w_0 by $w \mapsto f'(w_0)w + f(w_0) - f'(w_0)w_0 = f(w_0) + f'(w_0)(w - w_0)$.

infinitesimal variation δ around w_0 , the secant converges to the **tangent** of f at w_0 and the resulting slope defines the derivative of f at w_0 . The definition below formalizes this intuition.

Definition 2.3 (Derivative). The **derivative** of $f : \mathbb{R} \rightarrow \mathbb{R}$ at $w \in \mathbb{R}$ is defined as

$$f'(w) := \lim_{\delta \rightarrow 0} \frac{f(w + \delta) - f(w)}{\delta}, \quad (2.1)$$

provided that the limit exists. If $f'(w)$ is well-defined at a particular w , we say that the function f is **differentiable** at w . If f is differentiable at any $w \in \mathbb{R}$, we say that it is **differentiable everywhere** or differentiable for short.

If f is differentiable at a given w , then it is necessarily **continuous** at w as shown in the following proposition. Non-differentiability of a continuous function at a given point w is generally illustrated by a kink, as shown in Fig. 2.2.

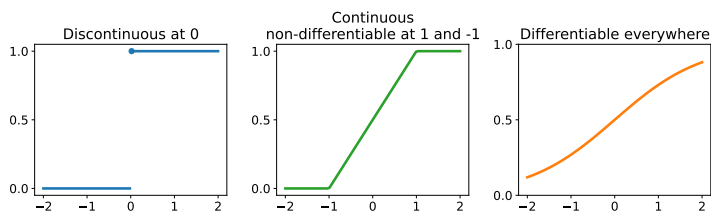


Figure 2.2: Illustration of discontinuity and non-differentiability. **Left.** A discontinuous function presents a jump in function values at a given point. **Center.** A continuous but non-differentiable everywhere function presents kinks at the points of non-differentiability. **Right.** A differentiable everywhere function is smooth.

Proposition 2.1 (Differentiability implies continuity). If $f : \mathbb{R} \rightarrow \mathbb{R}$ is differentiable at $w \in \mathbb{R}$, then it is continuous at $w \in \mathbb{R}$.

Proof. In little o notation, f is differentiable at w if there exists $f'(w) \in \mathbb{R}$, such that

$$f(w + \delta) = f(w) + f'(w)\delta + o(\delta) \text{ as } \delta \rightarrow 0.$$

Since $f'(w)\delta + o(\delta) = o(1)$ as $\delta \rightarrow 0$, f is continuous at w . \square

In addition to enabling the construction of a linear approximation of f in a neighborhood of w , since it is the slope of the tangent of f at w , the derivative f' informs us about the **monotonicity** of f around w . If $f'(w)$ is positive, the function is increasing around w . Conversely, if $f'(w)$ is negative, the function is decreasing. Such information can be used to develop iterative algorithms seeking to minimize f by computing iterates of the form $w_{t+1} = w_t - \gamma f'(w_t)$ for $\gamma > 0$, which move along descent directions of f around w_t .

For several elementary functions such as w^n , e^w , $\ln w$, $\cos w$ or $\sin w$, their derivatives can be obtained directly by applying the definition of the derivative in Eq. (2.1) as we now illustrate.

Example 2.1 (Derivative of power function). Consider $f(w) = w^n$

for $w \in \mathbb{R}$, $n \in \mathbb{N} \setminus \{0\}$. For any $\delta \in \mathbb{R}$, we have

$$\begin{aligned} \frac{f(w + \delta) - f(w)}{\delta} &= \frac{(w + \delta)^n - w^n}{\delta} \\ &= \frac{\sum_{k=0}^n \binom{n}{k} \delta^k w^{n-k} - w^n}{\delta} \\ &= \sum_{k=1}^n \binom{n}{k} \delta^{k-1} w^{n-k} \\ &= \binom{n}{1} w^{n-1} + \sum_{k=2}^n \binom{n}{k} \delta^{k-1} w^{n-k}, \end{aligned}$$

where, in the second line, we used the binomial theorem. Since $\binom{n}{1} = n$ and $\lim_{\delta \rightarrow 0} \sum_{k=2}^n \binom{n}{k} \delta^{k-1} w^{n-k} = 0$, we get $f'(w) = nw^{n-1}$.

Remark 2.2 (Functions on a subset \mathcal{U} of \mathbb{R}). For simplicity, we presented the definition of the derivative for a function defined on the whole set of real numbers \mathbb{R} . If a function $f : \mathcal{U} \rightarrow \mathbb{R}$ is defined on a subset $\mathcal{U} \subseteq \mathbb{R}$ of the real numbers, as it is the case for $f(w) = \sqrt{w}$ defined on $\mathcal{U} = \mathbb{R}_+$, the derivative of f at $w \in \mathcal{U}$ is defined by the limit in Eq. (2.1) provided that the function f is well defined on a neighborhood of w , that is, there exists $r > 0$ such that $w + \delta \in \mathcal{U}$ for any $|\delta| \leq r$. The function f is then said **differentiable everywhere** or differentiable for short if it is differentiable at any point w in the **interior** of \mathcal{U} , the set of points $w \in \mathcal{U}$ such that $\{w + \delta : |\delta| \leq r\} \subseteq \mathcal{U}$ for r sufficiently small. For points lying at the boundary of \mathcal{U} (such as a and b if $\mathcal{U} = [a, b]$), one may define the right and left derivatives of f at a and b , meaning that the limit is taken by approaching a from the right or b from the left.

2.1.2 Calculus rules

For a given $w \in \mathbb{R}$ and two functions $f : \mathbb{R} \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}$, the derivative of elementary operations on f and g such as their sums, products or compositions can easily be derived from the definition of the derivative, under appropriate conditions on the differentiability of f and g at w . For example, if the derivatives of f and g exist at w , then

the derivatives of their weighted sum or product exist, and satisfy the rules

$$\forall a, b \in \mathbb{R}, (af + bg)'(w) = af'(w) + bg'(w) \quad (\text{Linearity})$$

$$(fg)'(w) = f'(w)g(w) + f(w)g'(w), \quad (\text{Product rule})$$

where $(fg)(w) := f(w)g(w)$. The linearity can be verified directly from the linearity of the limits. For the product rule, in little o notation, we have, as $\delta \rightarrow 0$,

$$\begin{aligned} (fg)(w + \delta) &= (f(w) + f'(w)\delta + o(\delta))(g(w) + g'(w)\delta + o(\delta)) \\ &= f(w)g(w) + f'(w)g(w)\delta + f(w)g'(w)\delta + o(\delta), \end{aligned}$$

hence the result.

If the derivatives of g at w and of f at $g(w)$ exist, then the derivative of the composition $(f \circ g)(w) := f(g(w))$ at w exists and is given by

$$(f \circ g)'(w) = f'(g(w))g'(w). \quad (\text{Chain rule})$$

We prove this result more generally in Proposition 2.2. As seen in the sequel, the linearity and the product rule can be seen as byproducts of the chain rule, making the chain rule the cornerstone of differentiation.

Consider a function that can be expressed using sums, products or compositions of elementary functions, such as $f(w) := e^w \ln w + \cos w^2$. Its derivative can be computed by applying the aforementioned rules on the decomposition of f into elementary operations and functions.

Example 2.2 (Applying rules of differentiation). Consider $f(w) := e^w \ln w + \cos w^2$. The derivative of f at $w > 0$ can be computed step by step as follows, denoting $\text{sq}(w) := w^2$,

$$f'(w) = (\exp \cdot \ln)'(w) + (\cos \circ \text{sq})'(w) \quad (\text{Linearity})$$

$$(\exp \cdot \ln)'(w) = \exp'(w) \cdot \ln(w) + \exp(w) \cdot \ln'(w) \quad (\text{Product rule})$$

$$(\cos \circ \text{sq})'(w) = \cos'(\text{sq}(w)) \text{sq}'(w) \quad (\text{Chain rule})$$

$$\exp'(w) = \exp(w), \quad \ln'(w) = 1/w, \quad (\text{Elem. func.})$$

$$\text{sq}'(w) = 2w, \quad \cos'(w) = -\sin(w). \quad (\text{Elem. func.})$$

| We therefore obtain that $f'(w) = e^w \ln w + e^w/w - 2w \sin w^2$.

Such a process is purely mechanical and lends itself to an automated procedure, which is the main idea of automatic differentiation presented in Chapter 8.

2.1.3 Leibniz's notation

The notion of derivative was first introduced independently by Newton and Leibniz in the 18th century (Ball, 1960). The latter considered derivatives as the quotient of infinitesimal variations. Namely, denoting $u = f(w)$ a variable depending on w through f , Leibniz considered the derivative of f as the quotient

$$f' = \frac{du}{dw} \quad \text{with} \quad f'(w) = \left. \frac{du}{dw} \right|_w$$

where du and dw denote infinitesimal variations of u and w respectively and the symbol $|_w$ denotes the evaluation of the derivative at a given point w . This notation simplifies the statement of the chain rule first discovered by Leibniz (Rodriguez and Lopez Fernandez, 2010) as we have for $v = g(w)$ and $u = f(v)$

$$\frac{du}{dw} = \frac{du}{dv} \cdot \frac{dv}{dw}.$$

This hints that derivatives are multiplied when considering compositions. At evaluation, the chain rule in Leibniz notation recovers the formula presented above as

$$\left. \frac{du}{dw} \right|_w = \left. \frac{du}{dv} \right|_{g(w)} \left. \frac{dv}{dw} \right|_w = f'(g(w))g'(w) = (f \circ g)'(w).$$

The ability of Leibniz's notation to capture the chain rule as a mere product of quotients made it popular throughout the centuries, especially in mechanics (Ball, 1960). The rationale behind Leibniz's notation, the concept of "infinitesimal variations", was questioned by later mathematicians for its potential logical issues (Ball, 1960). The notation $f'(w)$ first introduced by Euler and further popularized by Lagrange (Cajori, 1993) has then taken over in numerous mathematical textbooks. The concept of infinitesimal variations has been rigorously defined by using

the set of hyperreal numbers. They extend the set of real numbers by considering each number as a sum of a non-infinitesimal part and an infinitesimal part (Hewitt, 1948). The formalism of infinitesimal variations further underlies the development of automatic differentiation algorithms through the concept of dual numbers.

2.2 Multivariate functions

2.2.1 Directional derivatives

Let us now consider a function $f : \mathbb{R}^P \rightarrow \mathbb{R}$ with multi-dimensional input $\mathbf{w} := (w_1, \dots, w_P) \in \mathbb{R}^P$. The most important example in machine learning is a function which, to the parameters $\mathbf{w} \in \mathbb{R}^P$ of a neural network, associates a loss value in \mathbb{R} . Variations of f need to be defined along specific directions, such as the variation $f(\mathbf{w} + \delta \mathbf{v}) - f(\mathbf{w})$ of f around $\mathbf{w} \in \mathbb{R}^P$ in the direction $\mathbf{v} \in \mathbb{R}^P$ by an amount $\delta > 0$. This consideration naturally leads to the definition of the directional derivative.

Definition 2.4 (Directional derivative). The **directional derivative** of f at \mathbf{w} in the **direction** \mathbf{v} is given by

$$\partial f(\mathbf{w})[\mathbf{v}] := \lim_{\delta \rightarrow 0} \frac{f(\mathbf{w} + \delta \mathbf{v}) - f(\mathbf{w})}{\delta},$$

provided that the limit exists.

We use the notation $[\mathbf{v}]$ to emphasize that, for a given input \mathbf{w} , we can see $\mathbf{v} \mapsto \partial f(\mathbf{w})[\mathbf{v}]$ as a function. This is essential to define the differentiability of f (Definition 2.6) at \mathbf{w} and to later define linear maps (Section 2.3).

One example of directional derivative consists in computing the derivative of a function f at \mathbf{w} in any of the canonical directions

$$\mathbf{e}_i := (0, \dots, 0, \underbrace{1}_i, 0, \dots, 0).$$

This allows us to define the notion of **partial derivatives**, denoted for $i \in [P]$

$$\partial_i f(\mathbf{w}) := \partial f(\mathbf{w})[\mathbf{e}_i] = \lim_{\delta \rightarrow 0} \frac{f(\mathbf{w} + \delta \mathbf{e}_i) - f(\mathbf{w})}{\delta}.$$

This is also denoted in Leibniz's notation as $\partial_i f(\mathbf{w}) = \frac{\partial f(\mathbf{w})}{\partial w_i}$ or $\partial_i f(\mathbf{w}) = \partial_{w_i} f(\mathbf{w})$. By moving along only the i^{th} coordinate of the function, the partial derivative is akin to differentiating the function $\omega_i \mapsto f(w_1, \dots, \omega_i, \dots, w_P)$ around ω_i , letting all other coordinates fixed at their values w_i .

2.2.2 Gradients

We now introduce the gradient vector, which gathers the partial derivatives. We first recall the definitions of linear map and linear form.

Definition 2.5 (Linear map, linear form). A function $l : \mathbb{R}^P \rightarrow \mathbb{R}^M$ is a **linear map** if for any $a_1, a_2 \in \mathbb{R}$, $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^D$,

$$l[a_1 \mathbf{v}_1 + a_2 \mathbf{v}_2] = a_1 l[\mathbf{v}_1] + a_2 l[\mathbf{v}_2].$$

A linear map with values in \mathbb{R} , $l : \mathbb{R}^P \rightarrow \mathbb{R}$, is called a **linear form**.

Linearity plays a crucial role in the differentiability of a function.

Definition 2.6 (Differentiability, single-output case). A function $f : \mathbb{R}^P \rightarrow \mathbb{R}$ is **differentiable** at $\mathbf{w} \in \mathbb{R}^P$ if its directional derivative is defined along any direction, is linear in any direction \mathbf{v} , and if

$$\lim_{\|\mathbf{v}\|_2 \rightarrow 0} \frac{|f(\mathbf{w} + \mathbf{v}) - f(\mathbf{w}) - \partial f(\mathbf{w})[\mathbf{v}]|}{\|\mathbf{v}\|_2} = 0.$$

We can now introduce the gradient.

Definition 2.7 (Gradient). The **gradient** of a differentiable function $f : \mathbb{R}^P \rightarrow \mathbb{R}$ at a point $\mathbf{w} \in \mathbb{R}^P$ is defined as the vector of partial derivatives

$$\nabla f(\mathbf{w}) := \begin{pmatrix} \partial_1 f(\mathbf{w}) \\ \vdots \\ \partial_P f(\mathbf{w}) \end{pmatrix} = \begin{pmatrix} \partial f(\mathbf{w})[\mathbf{e}_1] \\ \vdots \\ \partial f(\mathbf{w})[\mathbf{e}_P] \end{pmatrix} \in \mathbb{R}^P.$$

By linearity, the directional derivative of f at \mathbf{w} in the direction

$\mathbf{v} = \sum_{i=1}^P v_i \mathbf{e}_i$ is then given by

$$\partial f(\mathbf{w})[\mathbf{v}] = \sum_{i=1}^P v_i \partial f(\mathbf{w})[\mathbf{e}_i] = \langle \mathbf{v}, \nabla f(\mathbf{w}) \rangle \in \mathbb{R}.$$

Here, $\langle \cdot, \cdot \rangle$ denotes the inner product. We provide its definition in Euclidean spaces in Section 2.3.2.

In the definition above, the fact that the gradient can be used to compute the directional derivative is a mere consequence of the linearity of $\partial f(\mathbf{w})[\mathbf{v}]$ w.r.t. \mathbf{v} . However, in more abstract cases presented in later sections, the gradient is defined directly through this property.

As a simple example, any linear function of the form $f(\mathbf{w}) = \langle \mathbf{a}, \mathbf{w} \rangle = \sum_{i=1}^P a_i w_i$ is differentiable as we have $(\langle \mathbf{a}, \mathbf{w} + \mathbf{v} \rangle - \langle \mathbf{a}, \mathbf{w} \rangle - \langle \mathbf{a}, \mathbf{v} \rangle) / \|\mathbf{v}\|_2 = 0$ for any \mathbf{v} and in particular for $\|\mathbf{v}\| \rightarrow 0$. Moreover, its gradient is naturally given by $\nabla f(\mathbf{w}) = \mathbf{a}$.

More generally, to show that a function is differentiable and find its gradient, one approach is to approximate $f(\mathbf{w} + \mathbf{v})$ around $\mathbf{v} = \mathbf{0}$. If we can find a vector \mathbf{g} such that

$$f(\mathbf{w} + \mathbf{v}) = f(\mathbf{w}) + \langle \mathbf{g}, \mathbf{v} \rangle + o(\|\mathbf{v}\|_2),$$

then f is differentiable at \mathbf{w} , since $\langle \mathbf{g}, \cdot \rangle$ is linear. Moreover, \mathbf{g} is then the gradient of f at \mathbf{w} .

Remark 2.3 (Gateaux and Fréchet differentiability). Multiple definitions of differentiability exist. The one presented in Definition 2.6 is that of **Fréchet differentiable** functions. Alternatively, if $f : \mathbb{R}^P \rightarrow \mathbb{R}$ has well-defined directional derivatives along any directions then the function is **Gateaux differentiable**. Note that the existence of directional derivatives in any directions is not a sufficient condition for the function to be differentiable. In other words, any Fréchet differentiable function is Gateaux differentiable, but the converse is not true. As a counter-example, one can verify that the function $f(x_1, x_2) := x_1^3 / (x_1^2 + x_2^2)$ is Gateaux differentiable at 0 but not (Fréchet) differentiable at 0 (because the directional derivative at 0 is not linear).

Some authors also require Gateaux differentiable functions to

have linear directional derivatives along any direction. These are still not Fréchet differentiable functions. Indeed, the limit in Definition 2.6 is over any vectors tending to 0 (potentially in a pathological way), while directional derivatives look at such limits uniquely in terms of a single direction.

In the remainder of this chapter, all definitions of differentiability are in terms of Fréchet differentiability.

The next example illustrates how to compute the gradient of the logistic loss and validates its differentiability.

Example 2.3 (Gradient of logistic loss). Consider the logistic loss $\ell(\boldsymbol{\theta}, \mathbf{y}) := -\langle \mathbf{y}, \boldsymbol{\theta} \rangle + \log \sum_{i=1}^M e^{\theta_i}$, that measures the prediction error of the logits $\boldsymbol{\theta} \in \mathbb{R}^M$ w.r.t. the correct label $\mathbf{y} \in \{\mathbf{e}_1, \dots, \mathbf{e}_M\}$. Let us compute the gradient of this loss w.r.t. $\boldsymbol{\theta}$ for fixed \mathbf{y} , i.e., we want to compute the gradient of $f(\boldsymbol{\theta}) := \ell(\boldsymbol{\theta}, \mathbf{y})$. Let us decompose f as $f = l + \text{logsumexp}$ with $l(\boldsymbol{\theta}) := \langle -\mathbf{y}, \boldsymbol{\theta} \rangle$ and

$$\text{logsumexp}(\boldsymbol{\theta}) := \log \sum_{i=1}^M \exp(\theta_i),$$

the log-sum-exp function. The function l is linear so differentiable with gradient $\nabla l(\boldsymbol{\theta}) = -\mathbf{y}$. We therefore focus on logsumexp. Denoting $\exp(\boldsymbol{\theta}) = (\exp(\theta_1), \dots, \exp(\theta_M))$, using that $\exp(x) = 1 + x + o(x)$, $\log(1 + x) = x + o(x)$, and denoting \odot the elementwise product, we get

$$\begin{aligned} \text{logsumexp}(\boldsymbol{\theta} + \mathbf{v}) &= \log(\langle \exp(\boldsymbol{\theta} + \mathbf{v}), \mathbf{1} \rangle) \\ &= \log(\langle \exp(\boldsymbol{\theta}) \odot \exp(\mathbf{v}), \mathbf{1} \rangle) \\ &= \log(\langle \exp(\boldsymbol{\theta}) \odot (\mathbf{1} + \mathbf{v} + o(\|\mathbf{v}\|_2)), \mathbf{1} \rangle) \\ &= \log(\langle \exp(\boldsymbol{\theta}), \mathbf{1} \rangle + \langle \exp(\boldsymbol{\theta}), \mathbf{v} \rangle + o(\|\mathbf{v}\|_2)) \\ &= \log(\langle \exp(\boldsymbol{\theta}), \mathbf{1} \rangle) + \left\langle \frac{\exp(\boldsymbol{\theta})}{\langle \exp(\boldsymbol{\theta}), \mathbf{1} \rangle}, \mathbf{v} \right\rangle + o(\|\mathbf{v}\|_2). \end{aligned}$$

The above decomposition of $\text{logsumexp}(\boldsymbol{\theta} + \mathbf{v})$ shows that it is

differentiable, and that $\nabla \text{logsumexp}(\boldsymbol{\theta}) = \text{softargmax}(\boldsymbol{\theta})$, where

$$\text{softargmax}(\boldsymbol{\theta}) := \left(e^{\theta_1} / \left(\sum_{j=1}^M e^{\theta_j} \right), \dots, e^{\theta_M} / \left(\sum_{j=1}^M e^{\theta_j} \right) \right).$$

Overall, we get that $\nabla f(\boldsymbol{\theta}) = -\mathbf{y} + \text{softargmax}(\boldsymbol{\theta})$.

Linearity of gradients

The notion of differentiability for multi-input functions naturally inherits from the linearity of derivatives for single-input functions. For any $u_1, \dots, u_M \in \mathbb{R}$ and any multi-input functions f_1, \dots, f_M differentiable at \mathbf{w} , the function $u_1 f_1 + \dots + u_M f_M$ is differentiable at \mathbf{w} and its gradient is

$$\nabla(u_1 f_1 + \dots + u_M f_M)(\mathbf{w}) = u_1 \nabla f_1(\mathbf{w}) + \dots + u_M \nabla f_M(\mathbf{w}).$$

Why is the gradient useful?

When f is differentiable, we say that \mathbf{v} is an **ascent direction** of f from \mathbf{w} if

$$\langle \mathbf{v}, \nabla f(\mathbf{w}) \rangle > 0.$$

Conversely, we say that \mathbf{v} is a **descent direction** of f from \mathbf{w} if

$$\langle \mathbf{v}, \nabla f(\mathbf{w}) \rangle < 0.$$

Using this definition, the gradient leads to the **steepest ascent direction** of f from \mathbf{w} . To see why, we note that

$$\begin{aligned} \arg \max_{\mathbf{v} \in \mathbb{R}^P, \|\mathbf{v}\|_2 \leq 1} \langle \mathbf{v}, \nabla f(\mathbf{w}) \rangle &= \arg \max_{\mathbf{v} \in \mathbb{R}^P, \|\mathbf{v}\|_2 \leq 1} \partial f(\mathbf{w})[\mathbf{v}] \\ &= \nabla f(\mathbf{w}) / \|\nabla f(\mathbf{w})\|_2, \end{aligned}$$

where we assumed $\nabla f(\mathbf{w}) \neq \mathbf{0}$. The gradient $\nabla f(\mathbf{w})$ is orthogonal to the level set of the function (the set of points \mathbf{w} sharing the same value $f(\mathbf{w})$) and points towards higher values of f , as illustrated in Fig. 2.3. Conversely, the negative gradient $-\nabla f(\mathbf{w})$ points towards lower values of f . This observation motivates the development of optimization algorithms such as gradient descent. It is based on iteratively performing

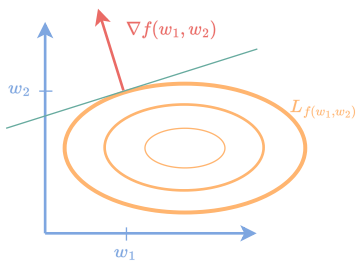


Figure 2.3: The gradient of a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ at (w_1, w_2) is the normal vector to the tangent space of the level set $L_{f(w_1, w_2)} = \{(w'_1, w'_2) : f(w'_1, w'_2) = f(w_1, w_2)\}$ and points towards points with higher function values.

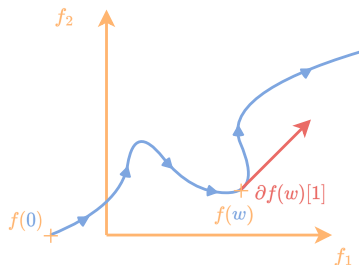


Figure 2.4: The directional derivative of a parametric curve $f : \mathbb{R} \rightarrow \mathbb{R}^2$ at w is the tangent to the curve at the point $f(w) \in \mathbb{R}^2$.

the update $\mathbf{w}_{t+1} := \mathbf{w}_t - \gamma \nabla f(\mathbf{w}_t)$, for $\gamma > 0$. It therefore seeks for a minimizer of f by moving along the **steepest descent direction** around \mathbf{w}_t given, up to a multiplicative factor, by $-\nabla f(\mathbf{w}_t)$. See also Definition 17.1 for more details.

2.2.3 Jacobians

Let us now consider a multi-output function $f : \mathbb{R}^P \rightarrow \mathbb{R}^M$ defined by $f(\mathbf{w}) := (f_1(\mathbf{w}), \dots, f_M(\mathbf{w}))$, where $f_j : \mathbb{R}^P \rightarrow \mathbb{R}$. A typical example in machine learning is a neural network. The notion of directional derivative can be extended to such function by defining it as the vector composed of the coordinate-wise directional derivatives:

$$\partial f(\mathbf{w})[\mathbf{v}] := \lim_{\delta \rightarrow 0} \frac{f(\mathbf{w} + \delta \mathbf{v}) - f(\mathbf{w})}{\delta} = \lim_{\delta \rightarrow 0} \begin{pmatrix} \frac{f_1(\mathbf{w} + \delta \mathbf{v}) - f_1(\mathbf{w})}{\delta} \\ \vdots \\ \frac{f_M(\mathbf{w} + \delta \mathbf{v}) - f_M(\mathbf{w})}{\delta} \end{pmatrix} \in \mathbb{R}^M,$$

where the limits (provided that they exist) are applied coordinate-wise. The directional derivative of f in the direction $\mathbf{v} \in \mathbb{R}^P$ is therefore the vector that gathers the directional derivative of each f_j , i.e., $\partial f(\mathbf{w})[\mathbf{v}] = (\partial f_j(\mathbf{w})[\mathbf{v}])_{j=1}^M$. In particular, we can define the **partial derivatives** of

f at \mathbf{w} as the vectors

$$\partial_i f(\mathbf{w}) := \partial f(\mathbf{w})[\mathbf{e}_i] = \begin{pmatrix} \partial_i f_1(\mathbf{w}) \\ \vdots \\ \partial_i f_M(\mathbf{w}) \end{pmatrix} \in \mathbb{R}^M.$$

As for the usual definition of the derivative, the directional derivative can provide a linear approximation of a function around a current input, as illustrated in Fig. 2.4 for a parametric curve $f : \mathbb{R} \rightarrow \mathbb{R}^2$.

Just as in the single-output case, differentiability is defined not only as the existence of directional derivatives in any direction but also by the linearity in the chosen direction.

Definition 2.8 (Differentiability, multi-output case). A function $f : \mathbb{R}^P \rightarrow \mathbb{R}^M$ is (Fréchet) **differentiable** at a point $\mathbf{w} \in \mathbb{R}^P$ if its directional derivative is defined along any direction, is linear in any direction, and,

$$\lim_{\|\mathbf{v}\|_2 \rightarrow 0} \frac{\|f(\mathbf{w} + \mathbf{v}) - f(\mathbf{w}) - \partial f(\mathbf{w})[\mathbf{v}]\|_2}{\|\mathbf{v}\|_2} = 0.$$

The partial derivatives of each coordinate's function are gathered in the **Jacobian matrix**.

Definition 2.9 (Jacobian). The **Jacobian** of a differentiable function $f : \mathbb{R}^P \rightarrow \mathbb{R}^M$ at \mathbf{w} is defined as the matrix gathering partial derivatives of each coordinate's function provided they exist,

$$\partial f(\mathbf{w}) := \begin{pmatrix} \partial_1 f_1(\mathbf{w}) & \dots & \partial_P f_1(\mathbf{w}) \\ \vdots & \ddots & \vdots \\ \partial_1 f_M(\mathbf{w}) & \dots & \partial_P f_M(\mathbf{w}) \end{pmatrix} \in \mathbb{R}^{M \times P}.$$

The Jacobian can be represented by stacking columns of partial derivatives or rows of gradients,

$$\partial f(\mathbf{w}) = \left(\partial_1 f(\mathbf{w}), \dots, \partial_P f(\mathbf{w}) \right) = \begin{pmatrix} \nabla f_1(\mathbf{w})^\top \\ \vdots \\ \nabla f_M(\mathbf{w})^\top \end{pmatrix} \in \mathbb{R}^{M \times P}.$$

By linearity, the directional derivative of f at \mathbf{w} along any input direction $\mathbf{v} = \sum_{i=1}^P v_i \mathbf{e}_i \in \mathbb{R}^P$ is then given by

$$\partial f(\mathbf{w})[\mathbf{v}] = \sum_{i=1}^P v_i \partial_i f(\mathbf{w}) = \partial f(\mathbf{w}) \mathbf{v} \in \mathbb{R}^M.$$

Notice that we use bold ∂ to indicate the Jacobian, seen as a matrix. The Jacobian matrix naturally generalizes the concepts of derivatives and gradients presented earlier. As for the single input case, to show that a function is differentiable, one approach is to approximate $f(\mathbf{w} + \mathbf{v})$ around $\mathbf{v} = \mathbf{0}$. If we find a linear map l such that

$$f(\mathbf{w} + \mathbf{v}) = f(\mathbf{w}) + l[\mathbf{v}] + o(\|\mathbf{v}\|_2),$$

then f is differentiable at \mathbf{w} . Moreover, if l is represented by matrix \mathbf{J} such that $l[\mathbf{v}] = \mathbf{J}\mathbf{v}$ then $\mathbf{J} = \partial f(\mathbf{w})$.

As a simple example, any linear function $f(\mathbf{w}) = \mathbf{A}\mathbf{w}$ for $\mathbf{A} \in \mathbb{R}^{M \times P}$ is differentiable, since all its coordinate-wise components are single-output linear functions, and the Jacobian of f at any \mathbf{w} is given by $\partial f(\mathbf{w}) = \mathbf{A}$.

Remark 2.4 (Special cases of the Jacobian). For single-output functions $f : \mathbb{R}^P \rightarrow \mathbb{R}$, i.e., $M = 1$, the Jacobian matrix reduces to a row vector identified as the **transpose of the gradient**,

$$\partial f(\mathbf{w}) = \nabla f(\mathbf{w})^\top \in \mathbb{R}^{1 \times P}.$$

For a single-input function $f : \mathbb{R} \rightarrow \mathbb{R}^M$, the Jacobian reduces to a single column vector of directional derivatives, denoted

$$\partial f(w) = f'(w) := \begin{pmatrix} f'_1(w) \\ \vdots \\ f'_M(w) \end{pmatrix} \in \mathbb{R}^{M \times 1}.$$

For a single-input single-output function $f : \mathbb{R} \rightarrow \mathbb{R}$, the Jacobian reduces to the derivative of f , i.e.,

$$\partial f(w) = f'(w) \in \mathbb{R}.$$

The next example illustrates the form of the Jacobian matrix for the element-wise application of a differentiable function σ , such as the softplus activation. In this case, the Jacobian takes a simple diagonal matrix form. As a consequence, the directional derivative associated with this function is simply given by an element-wise product: a full matrix-vector product is not needed, as would suggest Definition 2.9. We will revisit this point in Section 2.3.

Example 2.4 (Jacobian matrix of the softplus activation). Consider the element-wise application of the softplus defined for $\mathbf{w} \in \mathbb{R}^P$ by

$$f(\mathbf{w}) := \begin{pmatrix} \sigma(w_1) \\ \vdots \\ \sigma(w_P) \end{pmatrix} \in \mathbb{R}^P \quad \text{where} \quad \sigma(w) := \log(1 + e^w).$$

Since σ is differentiable, each coordinate of this function is differentiable and the overall function is differentiable. The j^{th} coordinate of f is independent of the i^{th} coordinate of \mathbf{w} for $i \neq j$, so $\partial_i f_j(\mathbf{w}) = 0$ for $i \neq j$. For $i = j$, the result boils down to the derivative of σ at w_j . That is, $\partial_j f_j(\mathbf{w}) = \sigma'(w_j)$, where $\sigma'(w) = e^w / (1 + e^w)$. The Jacobian of f is therefore a diagonal matrix

$$\partial f(\mathbf{w}) = \mathbf{diag}(\sigma'(w_1), \dots, \sigma'(w_P)) := \begin{pmatrix} \sigma'(w_1) & 0 & \dots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \sigma'(w_P) \end{pmatrix}.$$

Chain rule

Equipped with a generic definition of differentiability and the associated objects, gradients and Jacobians, we can now generalize the chain rule, previously introduced for single-input single-output functions.

Proposition 2.2 (Chain rule). Consider $f : \mathbb{R}^P \rightarrow \mathbb{R}^M$ and $g : \mathbb{R}^M \rightarrow \mathbb{R}^R$. If f is differentiable at $\mathbf{w} \in \mathbb{R}^P$ and g is differentiable at $f(\mathbf{w}) \in \mathbb{R}^M$, then the composition $g \circ f$ is differentiable

at $\mathbf{w} \in \mathbb{R}^P$ and its Jacobian is given by

$$\partial(g \circ f)(\mathbf{w}) = \partial g(f(\mathbf{w}))\partial f(\mathbf{w}).$$

Proof. We progressively approximate $g \circ f(\mathbf{w} + \mathbf{v})$ using the differentiability of f at \mathbf{w} and g at $f(\mathbf{w})$,

$$\begin{aligned} g(f(\mathbf{w} + \mathbf{v})) &= g(f(\mathbf{w}) + \partial f(\mathbf{w})\mathbf{v} + o(\|\mathbf{v}\|)) \\ &= g(f(\mathbf{w})) + \partial g(f(\mathbf{w}))\partial f(\mathbf{w})\mathbf{v} + o(\|\mathbf{v}\|). \end{aligned}$$

Hence, $g \circ f$ is differentiable at \mathbf{w} with Jacobian $\partial g(f(\mathbf{w}))\partial f(\mathbf{w})$. \square

Proposition 2.2 can be seen as the cornerstone of any derivative computations. For example, it can be used to rederive the linearity and product rules associated to the derivatives of single-input single-output functions.

When g is scalar-valued, combined with Remark 2.4, we obtain a simple expression for $\nabla(g \circ f)$.

Proposition 2.3 (Chain rule, scalar-valued case). Consider $f : \mathbb{R}^P \rightarrow \mathbb{R}^M$ and $g : \mathbb{R}^M \rightarrow \mathbb{R}$. The gradient of the composition is given by

$$\nabla(g \circ f)(\mathbf{w}) = \partial f(\mathbf{w})^\top \nabla g(f(\mathbf{w})).$$

This is a very useful identity in machine learning, as we often need to compose a vector-valued model function and a scalar-valued loss function. We illustrate this with linear regression below.

Example 2.5 (Linear regression). Consider a linear regression of N inputs $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$ onto N targets $y_1, \dots, y_N \in \mathbb{R}$, using a parameter vector $\mathbf{w} \in \mathbb{R}^D$. The loss is defined as the sum of squared residuals, $L(\mathbf{w}) := \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \sum_{i=1}^N (\langle \mathbf{x}_i, \mathbf{w} \rangle - y_i)^2$ where $\mathbf{X} := (\mathbf{x}_1, \dots, \mathbf{x}_N)^\top \in \mathbb{R}^{N \times D}$ and $\mathbf{y} := (y_1, \dots, y_N)^\top \in \mathbb{R}^N$.

The function L can be decomposed into a linear mapping $f(\mathbf{w}) := \mathbf{X}\mathbf{w}$ and a squared error $\ell(\mathbf{p}) := \|\mathbf{p} - \mathbf{y}\|_2^2$, so that $L = \ell \circ f$. We can then apply the chain rule in Proposition 2.3 to get

$$\nabla L(\mathbf{w}) = \partial f(\mathbf{w})^\top \nabla \ell(f(\mathbf{w}))$$

provided that f and ℓ are differentiable at \mathbf{w} and $f(\mathbf{w})$, respectively.

The function f is linear so differentiable with Jacobian $\partial f(\mathbf{w}) = \mathbf{X}$. On the other hand the partial derivatives of ℓ are given by $\partial_j \ell(\mathbf{p}) = 2(p_j - y_j)$ for $j \in \{1, \dots, N\}$. Therefore, ℓ is differentiable at any \mathbf{p} and its gradient is $\nabla \ell(\mathbf{p}) = 2(\mathbf{p} - \mathbf{y})$. By combining the two, we then get the gradient of L as

$$\nabla L(\mathbf{w}) = 2\mathbf{X}^\top (f(\mathbf{w}) - \mathbf{y}) = 2\mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y}).$$

2.3 Linear maps

The Jacobian matrix is useful as a representation of the partial derivatives. However, the core idea underlying the definition of differentiable functions, as well as their implementation in an autodiff framework, lies in the access to two key **linear maps**. These two maps encode infinitesimal variations along **input** or **output** directions and are referred to, respectively, as **Jacobian-vector product** (JVP) and **Vector-jacobian product** (VJP). This section formalizes these notions, in the context of Euclidean spaces.

2.3.1 The need for linear maps

So far, we have focused on functions $f: \mathbb{R}^P \rightarrow \mathbb{R}^M$, that take a vector as input and produce a vector as output. However, functions that use matrix or even tensor inputs/outputs are common place in neural networks. For example, consider the function of matrices of the form $f(\mathbf{W}) := \mathbf{W}\mathbf{x}$, where $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{W} \in \mathbb{R}^{M \times D}$. This function takes a matrix as input, not a vector. Of course, a matrix $\mathbf{W} \in \mathbb{R}^{M \times D}$ can always be “flattened” into a vector $\mathbf{w} \in \mathbb{R}^{MD}$, by stacking the columns of \mathbf{W} . We denote this operation by $\mathbf{w} = \text{vec}(\mathbf{W})$ and its inverse by $\mathbf{W} = \text{vec}^{-1}(\mathbf{w})$. We can then equivalently write $f(\mathbf{W})$ as $\tilde{f}(\mathbf{w}) = f(\text{vec}^{-1}(\mathbf{w})) = \text{vec}^{-1}(\mathbf{w})\mathbf{x}$, so that the previous framework applies. However, we will now see that this would be inefficient.

Indeed, the resulting Jacobian of \tilde{f} at any \mathbf{w} consists in a matrix of size $\mathbb{R}^{M \times MD}$, which, after some computations, can be observed to be mostly filled with zeros. Getting the directional derivative of f at $\mathbf{W} \in \mathbb{R}^{M \times D}$ in a direction $\mathbf{V} \in \mathbb{R}^{M \times D}$ would consist in (i) vectorizing

\mathbf{V} into $\mathbf{v} = \text{vec}(\mathbf{V})$, (ii) computing the matrix-vector product $\partial \tilde{f}(\mathbf{w})\mathbf{v}$ at a cost of $M^3 D^2$ computations (ignoring the fact that the Jacobian has many zero entries), (iii) re-shaping the result into a matrix.

On the other hand, since f is linear in its matrix input, we can infer that the directional derivative of f at any $\mathbf{W} \in \mathbb{R}^{M \times D}$ in any direction $\mathbf{V} \in \mathbb{R}^{M \times D}$ is simply given by the function itself applied on \mathbf{V} . Namely, we have $\partial f(\mathbf{W})[\mathbf{V}] = f(\mathbf{V}) = \mathbf{V}\mathbf{x}$, which is simple to implement and clearly only requires MD operations. Note that the cost would have been the same, had we ignored the non-zero entries of $\partial \tilde{f}(\mathbf{w})$. The point here is that by considering the operations associated to the differentiation of a function as linear maps rather than using the associated representation as a Jacobian matrix, we can efficiently exploit the underlying input or output space structure. To that end, we now recall the main abstractions necessary to extend the previous definitions in the context of Euclidean spaces.

2.3.2 Euclidean spaces

Linear spaces, a.k.a. **vector spaces**, are spaces equipped with and closed under an addition rule compatible with multiplication by a scalar (we limit ourselves to the field of reals). Namely, in a linear space \mathcal{E} , there exist operations $+$ and \cdot , such that for any $\mathbf{u}, \mathbf{v} \in \mathcal{E}$, and $a \in \mathbb{R}$, we have $\mathbf{u} + \mathbf{v} \in \mathcal{E}$ and $a \cdot \mathbf{u} \in \mathcal{E}$.

Euclidean spaces are linear spaces equipped with a basis $\mathbf{e}_1, \dots, \mathbf{e}_P \in \mathcal{E}$. Any element $\mathbf{v} \in \mathcal{E}$ can be decomposed as $\mathbf{v} = \sum_{i=1}^P v_i \mathbf{e}_i$ for some unique scalars $v_1, \dots, v_P \in \mathbb{R}$. A canonical example of Euclidean space is the set \mathbb{R}^P of all vectors of size P that we already covered. The set of matrices $\mathbb{R}^{P_1 \times P_2}$ of size $P_1 \times P_2$ is also naturally a Euclidean space generated by the set of canonical matrices $\mathbf{E}_{ij} \in \{0, 1\}^{P_1 \times P_2}$ for $i \in [P_1], j \in [P_2]$ filled with zero except at the $(i, j)^{\text{th}}$ entry filled with one. For example, $\mathbf{W} \in \mathbb{R}^{P_1 \times P_2}$ can be written $\mathbf{W} = \sum_{i,j=1}^{P_1, P_2} W_{ij} \mathbf{E}_{ij}$. Euclidean spaces are naturally equipped with a notion of inner product.

Definition 2.10 (Inner product). An **inner product** on a linear space \mathcal{E} is a function $\langle \cdot, \cdot \rangle : \mathcal{E} \times \mathcal{E} \rightarrow \mathbb{R}$ that is

- bilinear: $\mathbf{x} \mapsto \langle \mathbf{x}, \mathbf{w} \rangle$ and $\mathbf{y} \mapsto \langle \mathbf{v}, \mathbf{y} \rangle$ are linear for any $\mathbf{w}, \mathbf{v} \in \mathcal{E}$,
- symmetric: $\langle \mathbf{w}, \mathbf{v} \rangle = \langle \mathbf{v}, \mathbf{w} \rangle$ for any $\mathbf{w}, \mathbf{v} \in \mathcal{E}$,
- positive definite: $\langle \mathbf{w}, \mathbf{w} \rangle \geq 0$ for any $\mathbf{w} \in \mathcal{E}$, and $\langle \mathbf{w}, \mathbf{w} \rangle = 0$ if and only if $\mathbf{w} = 0$.

An inner product defines a norm $\|\mathbf{w}\| := \sqrt{\langle \mathbf{w}, \mathbf{w} \rangle}$.

The norm induced by an inner product defines a distance $\|\mathbf{w} - \mathbf{v}\|$ between $\mathbf{w}, \mathbf{v} \in \mathcal{E}$, and therefore a notion of convergence.

For vectors, where $\mathcal{E} = \mathbb{R}^P$, the inner product is the usual one $\langle \mathbf{w}, \mathbf{v} \rangle = \sum_{i=1}^P w_i v_i$. For matrices, where $\mathcal{E} = \mathbb{R}^{P_1 \times P_2}$, the inner product is the so-called Frobenius inner product. It is defined for any $\mathbf{W}, \mathbf{V} \in \mathbb{R}^{P_1 \times P_2}$ by

$$\langle \mathbf{W}, \mathbf{V} \rangle := \langle \text{vec}(\mathbf{W}), \text{vec}(\mathbf{V}) \rangle = \sum_{i,j=1}^{P_1, P_2} W_{ij} V_{ij} = \text{tr}(\mathbf{W}^\top \mathbf{V}),$$

where $\text{tr}(\mathbf{Z}) := \sum_{i=1}^P Z_{ii}$ is the trace operator defined for square matrices $\mathbf{Z} \in \mathbb{R}^{P \times P}$. For tensors of order R , which generalize matrices to $\mathcal{E} = \mathbb{R}^{P_1 \times \dots \times P_R}$, the inner product is defined similarly for $\mathbf{W}, \mathbf{V} \in \mathbb{R}^{P_1 \times \dots \times P_R}$ by

$$\langle \mathbf{W}, \mathbf{V} \rangle := \langle \text{vec}(\mathbf{W}), \text{vec}(\mathbf{V}) \rangle = \sum_{i_1, \dots, i_R=1}^{P_1, \dots, P_R} \mathbf{W}_{i_1 \dots i_R} \mathbf{V}_{i_1 \dots i_R},$$

where $\mathbf{W}_{i_1 \dots i_R}$ is the $(i_1, \dots, i_R)^{\text{th}}$ entry of \mathbf{W} .

2.3.3 Linear maps and their adjoints

The notion of linear map in Definition 2.5 naturally extends to Euclidean spaces. Namely, a function $l : \mathcal{E} \rightarrow \mathcal{F}$ from a Euclidean space \mathcal{E} onto a Euclidean space \mathcal{F} is a **linear map** if for any $\mathbf{w}, \mathbf{v} \in \mathcal{E}$ and $a, b \in \mathbb{R}$, we have $l[a\mathbf{w} + b\mathbf{v}] = a \cdot l[\mathbf{w}] + b \cdot l[\mathbf{v}]$. When $\mathcal{E} = \mathbb{R}^P$ and $\mathcal{F} = \mathbb{R}^M$, there always exists a matrix $\mathbf{A} \in \mathbb{R}^{M \times P}$ such that $l[\mathbf{v}] = \mathbf{A}\mathbf{v}$. Therefore, we can think of \mathbf{A} as the “materialization” as a matrix of l .

Example 2.6 (Linear map). Consider the linear map $l[\mathbf{v}] := (\mathbf{a}\mathbf{b}^\top)\mathbf{v}$, where $\mathbf{a} \in \mathbb{R}^M$, $\mathbf{b} \in \mathbb{R}^P$ and $\mathbf{v} \in \mathbb{R}^P$. This is a function from \mathbb{R}^P to \mathbb{R}^M . We can always materialize the linear map as a matrix $\mathbf{A} := \mathbf{a}\mathbf{b}^\top \in \mathbb{R}^{M \times P}$ and write $l[\mathbf{v}] = \mathbf{A}\mathbf{v}$. However, applying a linear map on a vector \mathbf{v} often does not require materializing the corresponding matrix. Following the previous example, we can simply write $l[\mathbf{v}] = (\mathbf{b}^\top \mathbf{v})\mathbf{a}$, which only requires an inner product and an element-wise multiplication. This is more efficient than materializing \mathbf{A} then computing $\mathbf{A}\mathbf{v}$, which requires an outer product and a matrix-vector multiplication.

We can define the adjoint operator of a linear map.

Definition 2.11 (Adjoint operator). Given two Euclidean spaces \mathcal{E} and \mathcal{F} equipped with inner products $\langle \cdot, \cdot \rangle_{\mathcal{E}}$ and $\langle \cdot, \cdot \rangle_{\mathcal{F}}$, the **adjoint** of a linear map $l : \mathcal{E} \rightarrow \mathcal{F}$ is the unique linear map $l^* : \mathcal{F} \rightarrow \mathcal{E}$ such that for any $\mathbf{v} \in \mathcal{E}$ and $\mathbf{u} \in \mathcal{F}$,

$$\langle l[\mathbf{v}], \mathbf{u} \rangle_{\mathcal{F}} = \langle \mathbf{v}, l^*[\mathbf{u}] \rangle_{\mathcal{E}}.$$

The adjoint can be thought as the counterpart of the matrix transpose for linear maps. When $l[\mathbf{v}] = \mathbf{A}\mathbf{v}$, we have $l^*[\mathbf{u}] = \mathbf{A}^\top \mathbf{u}$ since

$$\langle l[\mathbf{v}], \mathbf{u} \rangle_{\mathcal{F}} = \langle \mathbf{A}\mathbf{v}, \mathbf{u} \rangle_{\mathcal{F}} = \langle \mathbf{v}, \mathbf{A}^\top \mathbf{u} \rangle_{\mathcal{E}} = \langle \mathbf{v}, l^*[\mathbf{u}] \rangle_{\mathcal{E}}.$$

Example 2.7 (Adjoint linear map). Using the linear map $l[\mathbf{v}]$ from the previous example, we have for all $\mathbf{u} \in \mathbb{R}^M$ and $\mathbf{v} \in \mathbb{R}^P$,

$$\langle \mathbf{a}\mathbf{b}^\top \mathbf{v}, \mathbf{u} \rangle = \langle \mathbf{v}, \mathbf{b}\mathbf{a}^\top \mathbf{u} \rangle.$$

Therefore, the adjoint linear map is $l^*[\mathbf{u}] = (\mathbf{b}\mathbf{a}^\top)\mathbf{u}$. This is a function from \mathbb{R}^M to \mathbb{R}^P . It can be materialized as the matrix $\mathbf{A}^\top = \mathbf{b}\mathbf{a}^\top \in \mathbb{R}^{P \times M}$. Applying $l^*[\mathbf{u}]$ can be done efficiently as $l^*[\mathbf{u}] = (\mathbf{a}^\top \mathbf{u})\mathbf{b}$.

2.3.4 Jacobian-vector products

We now define the directional derivative using linear maps, leading to the notion of Jacobian-vector product (JVP). This can be used to

facilitate the treatment of functions on tensors or for further extensions to infinite-dimensional spaces. In the following, \mathcal{E} and \mathcal{F} denote two Euclidean spaces equipped with inner products $\langle \cdot, \cdot \rangle_{\mathcal{E}}$ and $\langle \cdot, \cdot \rangle_{\mathcal{F}}$. We start by defining differentiability in general Euclidean spaces.

Definition 2.12 (Differentiability in Euclidean spaces). A function $f : \mathcal{E} \rightarrow \mathcal{F}$ is **differentiable** at a point $\mathbf{w} \in \mathcal{E}$ if the **directional derivative** along $\mathbf{v} \in \mathcal{E}$

$$\partial f(\mathbf{w})[\mathbf{v}] := \lim_{\delta \rightarrow 0} \frac{f(\mathbf{w} + \delta \mathbf{v}) - f(\mathbf{w})}{\delta}$$

is well-defined for any $\mathbf{v} \in \mathcal{E}$, linear in \mathbf{v} and if

$$\lim_{\|\mathbf{v}\|_{\mathcal{F}} \rightarrow 0} \frac{\|f(\mathbf{w} + \mathbf{v}) - f(\mathbf{w}) - \partial f(\mathbf{w})[\mathbf{v}]\|_{\mathcal{F}}}{\|\mathbf{v}\|_{\mathcal{F}}} = 0.$$

We can now formally define the Jacobian-vector product.

Definition 2.13 (Jacobian-vector product). For a differentiable function $f : \mathcal{E} \rightarrow \mathcal{F}$, the **linear map** $\partial f(\mathbf{w}) : \mathcal{E} \rightarrow \mathcal{F}$, mapping \mathbf{v} to $\partial f(\mathbf{w})[\mathbf{v}]$, is called the **Jacobian-vector product** (JVP). From this perspective, the function ∂f is a function from \mathcal{E} to a linear map from \mathcal{E} to \mathcal{F} . That is, we have

$$\partial f : \mathcal{E} \rightarrow (\mathcal{E} \rightarrow \mathcal{F}).$$

We emphasize again that the directional derivative $\partial f(\mathbf{w})[\mathbf{v}] \in \mathcal{F}$ is a value, while the JVP $\mathbf{v} \mapsto \partial f(\mathbf{w})[\mathbf{v}]$ is a function. Strictly speaking, \mathbf{v} can belong to any Euclidean space \mathcal{E} and does not need to be limited to a vector, as the JVP acronym would suggest. We adopt the name JVP, as it is now standard.

Recovering the gradient

Previously, we saw that for differentiable functions with vector input and scalar output, the directional derivative is equal to the inner product between the direction and the gradient. The same applies when considering differentiable functions from a Euclidean space with single

outputs, except that the gradient is now an element of the input space and the inner product is the one associated with the input space.

Proposition 2.4 (Gradient). If a function $f : \mathcal{E} \rightarrow \mathbb{R}$ is differentiable at $\mathbf{w} \in \mathcal{E}$, then there exists $\nabla f(\mathbf{w}) \in \mathcal{E}$, called the **gradient** of f at \mathbf{w} such that the directional derivative of f at \mathbf{w} along any input direction $\mathbf{v} \in \mathcal{E}$ is given by

$$\partial f(\mathbf{w})[\mathbf{v}] = \langle \nabla f(\mathbf{w}), \mathbf{v} \rangle_{\mathcal{E}}.$$

In Euclidean spaces, the existence of the gradient can simply be shown by decomposing the partial derivative along a basis of \mathcal{E} . Such a definition generalizes to infinite-dimensional (e.g., Hilbert spaces) spaces as discussed in Section 2.3.9.

2.3.5 Vector-Jacobian products

Consider a function $f : \mathbb{R}^P \rightarrow \mathbb{R}^M$. Instead of variations of f along an **input** direction $\mathbf{v} \in \mathbb{R}^P$, we may also consider the variations of f along an **output** direction $\mathbf{u} \in \mathbb{R}^M$, namely, computing the gradient $\nabla \langle \mathbf{u}, f \rangle(\mathbf{w})$ of the scalar-valued function

$$\langle \mathbf{u}, f \rangle(\mathbf{w}) := \langle \mathbf{u}, f(\mathbf{w}) \rangle \in \mathbb{R}.$$

Equivalently, we may compute the gradients $\nabla f_j(\mathbf{w})$ of each coordinate function $f_j := \langle \mathbf{e}_j, f \rangle$ at \mathbf{w} , where \mathbf{e}_j is the j^{th} canonical vector in \mathbb{R}^M . The infinitesimal variations of f at \mathbf{w} along any output direction $\mathbf{u} = \sum_{j=1}^M u_j \mathbf{e}_j \in \mathbb{R}^M$ are given by

$$\nabla \langle \mathbf{u}, f \rangle(\mathbf{w}) = \sum_{j=1}^M u_j \nabla f_j(\mathbf{w}) = \partial f(\mathbf{w})^\top \mathbf{u} \in \mathbb{R}^P,$$

where $\partial f(\mathbf{w})^\top \in \mathbb{R}^{P \times M}$ is the Jacobian's transpose. Using the definition of derivative as a limit, we may also write for $i \in [P]$

$$\nabla_i \langle \mathbf{u}, f \rangle(\mathbf{w}) = [\partial f(\mathbf{w})^\top \mathbf{u}]_i = \lim_{\delta \rightarrow 0} \frac{\langle \mathbf{u}, f(\mathbf{w} + \delta \mathbf{e}_i) - f(\mathbf{w}) \rangle}{\delta},$$

where \mathbf{e}_i is the i^{th} canonical vector in \mathbb{R}^P .

For generic Euclidean spaces \mathcal{E} and \mathcal{F} , the counterpart of the transpose is the adjoint operator, leading to the notion of vector-Jacobian product.

Proposition 2.5 (Vector-Jacobian product). If a function $f : \mathcal{E} \rightarrow \mathcal{F}$ is differentiable at $\mathbf{w} \in \mathcal{E}$, then its infinitesimal variation along an **output** direction $\mathbf{u} \in \mathcal{F}$ is given by the **adjoint map** $\partial f(\mathbf{w})^* : \mathcal{F} \rightarrow \mathcal{E}$ of the JVP, called the **vector-Jacobian product** (VJP). It satisfies

$$\nabla \langle \mathbf{u}, f \rangle_{\mathcal{F}}(\mathbf{w}) = \partial f(\mathbf{w})^*[\mathbf{u}],$$

where we denoted $\langle \mathbf{u}, f \rangle_{\mathcal{F}}(\mathbf{w}) := \langle \mathbf{u}, f(\mathbf{w}) \rangle_{\mathcal{F}}$. The function $\partial f(\cdot)^*$ is a function from \mathcal{E} to a linear map from \mathcal{F} to \mathcal{E} . That is, we have

$$\partial f(\cdot)^* : \mathcal{E} \rightarrow (\mathcal{F} \rightarrow \mathcal{E}).$$

Proof. The chain rule presented in Proposition 2.2 naturally generalizes to Euclidean spaces (see Proposition 2.6). Since $\langle \mathbf{u}, \cdot \rangle_{\mathcal{F}}$ is linear, its directional derivative is itself. Therefore, the directional derivative of $\langle \mathbf{u}, f \rangle_{\mathcal{F}}$ is

$$\begin{aligned} \partial(\langle \mathbf{u}, f \rangle_{\mathcal{F}})(\mathbf{w})[\mathbf{v}] &= \langle \mathbf{u}, \partial f(\mathbf{w})[\mathbf{v}] \rangle_{\mathcal{F}} \\ &= \langle \partial f(\mathbf{w})^*[\mathbf{u}], \mathbf{v} \rangle_{\mathcal{E}}. \end{aligned}$$

As this is true for any $\mathbf{v} \in \mathcal{E}$, $\partial f(\mathbf{w})^*[\mathbf{u}]$ is the gradient of $\langle \mathbf{u}, f \rangle_{\mathcal{F}}$ per Proposition 2.4. \square

We illustrate the JVP and VJP linear maps in Fig. 2.5.

2.3.6 Chain rule using linear maps

The chain rule presented before in terms of Jacobian matrices can readily be formulated to take advantage of the implementations of the JVP and VJP as linear maps.

Proposition 2.6 (Chain rule, general case). Consider $f : \mathcal{E} \rightarrow \mathcal{F}$ and $g : \mathcal{F} \rightarrow \mathcal{G}$, where \mathcal{E} , \mathcal{F} and \mathcal{G} are Euclidean spaces. If f is differentiable at $\mathbf{w} \in \mathcal{E}$ and g is differentiable at $f(\mathbf{w}) \in \mathcal{F}$, then

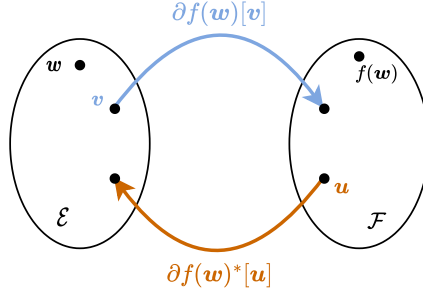


Figure 2.5: Jacobian-vector product (JVP) $v \mapsto \partial f(\mathbf{w})[v]$ and vector-Jacobian product (VJP) $u \mapsto \partial f(\mathbf{w})^*[u]$, seen as linear maps.

the composition $g \circ f$ is differentiable at $\mathbf{w} \in \mathcal{E}$. Its JVP is given for all $\mathbf{v} \in \mathcal{E}$ by

$$\partial(g \circ f)(\mathbf{w})[\mathbf{v}] = \partial g(f(\mathbf{w}))[\partial f(\mathbf{w})[\mathbf{v}]]$$

and its VJP is given for all $\mathbf{u} \in \mathcal{G}$ by

$$\partial(g \circ f)(\mathbf{w})^*[\mathbf{u}] = \partial f(\mathbf{w})^*[\partial g(f(\mathbf{w}))^*[\mathbf{u}]].$$

The proof follows the one of Proposition 2.2. This is illustrated in Fig. 2.6. When the last function is scalar-valued, which is often the case in machine learning, we obtain the following simplified result.

Proposition 2.7 (Chain rule, scalar case). Consider $f : \mathcal{E} \rightarrow \mathcal{F}$ and $g : \mathcal{F} \rightarrow \mathbb{R}$, the gradient of the composition is given by

$$\nabla(g \circ f)(\mathbf{w}) = \partial f(\mathbf{w})^*[\nabla g(f(\mathbf{w}))].$$

2.3.7 Functions of multiple inputs (fan-in)

Oftentimes, the inputs of a function do not belong to only one Euclidean space but to a product of them. An example is $f(\mathbf{x}, \mathbf{W}) := \mathbf{W}\mathbf{x}$, which is defined on $\mathcal{E} := \mathbb{R}^D \times \mathbb{R}^{M \times D}$. In such a case, it is convenient to generalize the notion of partial derivatives to handle blocks of inputs.

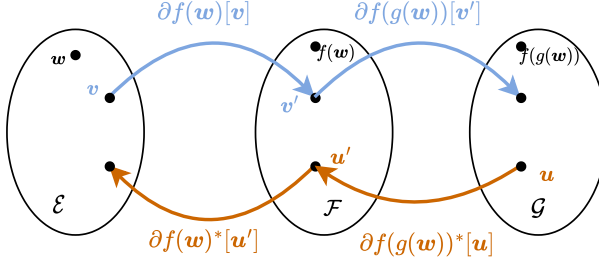


Figure 2.6: Chain rule using JVP and VJP linear maps.

Consider a function $f(\mathbf{w}_1, \dots, \mathbf{w}_S)$ defined on $\mathcal{E} := \mathcal{E}_1 \times \dots \times \mathcal{E}_S$, where $\mathbf{w}_i \in \mathcal{E}_i$. We denote the partial derivative with respect to the i^{th} input \mathbf{w}_i along $\mathbf{v}_i \in \mathcal{E}_i$ as $\partial_i f(\mathbf{w}_1, \dots, \mathbf{w}_S)[\mathbf{v}_i]$. Equipped with this notation, we can analyze how JVPs or VJPs are decomposed along several inputs.

Proposition 2.8 (Multiple inputs). Consider a differentiable function of the form $f(\mathbf{w}) = f(\mathbf{w}_1, \dots, \mathbf{w}_S)$ with signature $f: \mathcal{E} \rightarrow \mathcal{F}$, where $\mathbf{w} := (\mathbf{w}_1, \dots, \mathbf{w}_S) \in \mathcal{E}$ and $\mathcal{E} := \mathcal{E}_1 \times \dots \times \mathcal{E}_S$. Then the JVP with the input direction $\mathbf{v} = (\mathbf{v}_1, \dots, \mathbf{v}_S) \in \mathcal{E}$ is given by

$$\begin{aligned} \partial f(\mathbf{w})[\mathbf{v}] &= \partial f(\mathbf{w}_1, \dots, \mathbf{w}_S)[\mathbf{v}_1, \dots, \mathbf{v}_S] \in \mathcal{F} \\ &= \sum_{i=1}^S \partial_i f(\mathbf{w}_1, \dots, \mathbf{w}_S)[\mathbf{v}_i]. \end{aligned}$$

The VJP with the output direction $\mathbf{u} \in \mathcal{F}$ is given by

$$\begin{aligned} \partial f(\mathbf{w})^*[\mathbf{u}] &= \partial f(\mathbf{w}_1, \dots, \mathbf{w}_S)^*[\mathbf{u}] \in \mathcal{E} \\ &= (\partial_1 f(\mathbf{w}_1, \dots, \mathbf{w}_S)^*[\mathbf{u}], \dots, \partial_S f(\mathbf{w}_1, \dots, \mathbf{w}_S)^*[\mathbf{u}]). \end{aligned}$$

Example 2.8 (Matrix-vector product). Consider $f(\mathbf{x}, \mathbf{W}) := \mathbf{W}\mathbf{x}$, where $\mathbf{W} \in \mathbb{R}^{M \times D}$ and $\mathbf{x} \in \mathbb{R}^D$. This corresponds to setting $\mathcal{E} := \mathcal{E}_1 \times \mathcal{E}_2 := \mathbb{R}^D \times \mathbb{R}^{M \times D}$ and $\mathcal{F} := \mathbb{R}^M$. For the JVP, letting

$\mathbf{v} \in \mathcal{E}_1$ and $\mathbf{V} \in \mathcal{E}_2$, we obtain

$$\partial f(\mathbf{x}, \mathbf{W})[\mathbf{v}, \mathbf{V}] = \mathbf{W}\mathbf{v} + \mathbf{V}\mathbf{x} \in \mathcal{F}.$$

We can also access the individual JVPs as

$$\partial_1 f(\mathbf{x}, \mathbf{W})[\mathbf{v}] = \mathbf{W}\mathbf{v} \in \mathcal{F},$$

$$\partial_2 f(\mathbf{x}, \mathbf{W})[\mathbf{V}] = \mathbf{V}\mathbf{x} \in \mathcal{F}.$$

For the VJP, letting $\mathbf{u} \in \mathcal{F}$, we obtain

$$\partial f(\mathbf{x}, \mathbf{W})^*[\mathbf{u}] = (\mathbf{W}^\top \mathbf{u}, \mathbf{u}\mathbf{x}^\top) \in \mathcal{E}.$$

We can access the individual VJPs by

$$\partial_1 f(\mathbf{x}, \mathbf{W})^*[\mathbf{u}] = \mathbf{W}^\top \mathbf{u} \in \mathcal{E}_1,$$

$$\partial_2 f(\mathbf{x}, \mathbf{W})^*[\mathbf{u}] = \mathbf{u}\mathbf{x}^\top \in \mathcal{E}_2.$$

Remark 2.5 (Nested inputs). It is sometimes convenient to group inputs into meaningful parts. For instance, if the input is naturally broken down into two parts $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$, where \mathbf{x}_1 is a text part and \mathbf{x}_2 is an image part, and the network parameters are naturally grouped into three layers $\mathbf{w} = (\mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3)$, we can write $f(\mathbf{x}, \mathbf{w}) = f((\mathbf{x}_1, \mathbf{x}_2), (\mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3))$. This is mostly a convenience and we can again reduce it to a function of a single input, thanks to the linear map perspective in Euclidean spaces.

Remark 2.6 (Hiding away inputs). It will often be convenient to ignore inputs when differentiating. We use the semicolon for this purpose. For instance, a function of the form $L(\mathbf{w}; \mathbf{x}, \mathbf{y})$ (notice the semicolon) has signature $L: \mathcal{W} \rightarrow \mathbb{R}$ because we treat \mathbf{x} and \mathbf{y} as constants. Therefore, the gradient is $\nabla L(\mathbf{w}; \mathbf{x}, \mathbf{y}) \in \mathcal{W}$. On the other hand, the function $L(\mathbf{w}, \mathbf{x}, \mathbf{y})$ (notice the comma) has signature $L: \mathcal{W} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ so its gradient is $\nabla L(\mathbf{w}, \mathbf{x}, \mathbf{y}) \in \mathcal{W} \times \mathcal{X} \times \mathcal{Y}$. If we need to access partial gradients, we use indexing, e.g., $\nabla_1 L(\mathbf{w}, \mathbf{x}, \mathbf{y}) \in \mathcal{W}$ or $\nabla_{\mathbf{w}} L(\mathbf{w}, \mathbf{x}, \mathbf{y}) \in \mathcal{W}$ when there is no ambiguity.

2.3.8 Functions with multiple outputs (fan-out)

Similarly, it is often convenient to deal with functions that have multiple outputs.

Proposition 2.9 (Multiple outputs). Consider a differentiable function of the form $f(\mathbf{w}) := (f_1(\mathbf{w}), \dots, f_T(\mathbf{w}))$, with signatures $f: \mathcal{E} \rightarrow \mathcal{F}$ and $f_i: \mathcal{E} \rightarrow \mathcal{F}_i$, where $\mathcal{F} := \mathcal{F}_1 \times \dots \times \mathcal{F}_T$. Then the JVP with the input direction $\mathbf{v} \in \mathcal{E}$ is given by

$$\partial f(\mathbf{w})[\mathbf{v}] = (\partial f_1(\mathbf{w})[\mathbf{v}], \dots, \partial f_T(\mathbf{w})[\mathbf{v}]) \in \mathcal{F}.$$

The VJP with the output direction $\mathbf{u} = (\mathbf{u}_1, \dots, \mathbf{u}_T) \in \mathcal{F}$ is

$$\begin{aligned} \partial f(\mathbf{w})^*[\mathbf{u}] &= \partial f(\mathbf{w})^*[\mathbf{u}_1, \dots, \mathbf{u}_T] \in \mathcal{E} \\ &= \sum_{i=1}^T \partial f_i(\mathbf{w})^*[\mathbf{u}_i]. \end{aligned}$$

Combined with the chain rule, we obtain that the Jacobian of

$$h(\mathbf{w}) := g(f(\mathbf{w})) = g(f_1(\mathbf{w}), \dots, f_T(\mathbf{w}))$$

is $\partial h(\mathbf{w}) = \sum_{i=1}^T \partial_i g(f(\mathbf{w})) \circ \partial f_i(\mathbf{w})$ and therefore the JVP is

$$\partial h(\mathbf{w})[\mathbf{v}] = \sum_{i=1}^T \partial_i g(f(\mathbf{w}))[\partial f_i(\mathbf{w})[\mathbf{v}]].$$

2.3.9 Extensions to non-Euclidean linear spaces

So far, we focused on Euclidean spaces, i.e., linear spaces with a finite basis. However, the notions studied earlier can be generalized to more generic spaces.

For example, **directional derivatives** (see Definition 2.12) can be defined in any linear space equipped with a norm and complete with respect to this norm. Such spaces are called **Banach spaces**. Completeness is a technical assumption that requires that any Cauchy sequence converges (a Cauchy sequence is a sequence whose elements become arbitrarily close to each other as the sequence progresses). A function $f: \mathcal{E} \rightarrow \mathcal{F}$ defined from a Banach space \mathcal{E} onto a Banach space

\mathcal{F} is then called **Gateaux differentiable** if its directional derivative is defined along any direction (where limits are defined w.r.t. the norm in \mathcal{F}). Some authors also require the directional derivative to be linear to define a Gateaux differentiable function.

Fréchet differentiability can also naturally be generalized to Banach spaces. The only difference is that, in generic Banach spaces, the linear map l satisfying Definition 2.12 must be continuous, i.e., there must exist $C > 0$, such that $l[v] \leq C\|v\|$, where $\|\cdot\|$ is the norm in the Banach space \mathcal{E} .

The definitions of gradient and VJPs require in addition a notion of inner product. They can be defined in **Hilbert spaces**, that is, linear spaces equipped with an inner product and complete with respect to the norm induced by the inner product (they could also be defined in a Banach space by considering operations in the dual space, see, e.g. (Clarke *et al.*, 2008)). The existence of the gradient is ensured by **Riesz's representation theorem** which states that any continuous linear form in a Hilbert space can be represented by the inner product with a vector. Since for a differentiable function $f : \mathcal{E} \rightarrow \mathbb{R}$, the JVP $\partial f(\mathbf{w}) : \mathcal{E} \rightarrow \mathbb{R}$ is a linear form, Riesz's representation theorem ensures the existence of the gradient as the element $g \in \mathcal{E}$ such that $\partial f(\mathbf{w})\mathbf{v} = \langle g, \mathbf{v} \rangle$ for any $\mathbf{v} \in \mathcal{E}$. The VJP is also well-defined as the adjoint of the JVP w.r.t. the inner product of the Hilbert space.

As an example, the space of squared integrable functions on \mathbb{R} is a Hilbert space equipped with the inner product $\langle a, b \rangle := \int a(x)b(x)dx$. Here, we cannot find a finite number of functions that can express all possible functions on \mathbb{R} . Therefore, this space is not a mere Euclidean space. Nevertheless, we can consider functions on this Hilbert space (called **functionals** to distinguish them from the elements of the space). The associated directional derivatives and gradients, can be defined and are called respectively, **functional derivative** and **functional gradient**, see, e.g., Frigyik *et al.* (2008) and references therein.

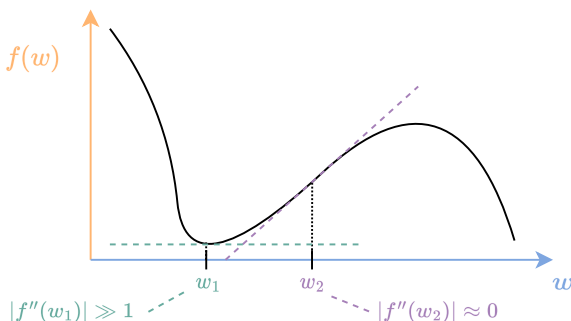


Figure 2.7: Points at which the second derivative is small are points along which the function is well approximated by its tangent line. On the other hand, points with large second derivative tend to be badly approximated by the tangent line.

2.4 Second-order differentiation

2.4.1 Second derivatives

For a single-input, single-output differentiable function $f : \mathbb{R} \rightarrow \mathbb{R}$, its derivative at any point is itself a function $f' : \mathbb{R} \rightarrow \mathbb{R}$. We may then consider the derivative of the derivative at any point: the **second derivative**.

Definition 2.14 (Second derivative). The **second derivative** $f^{(2)}(w)$ of a differentiable function $f : \mathbb{R} \rightarrow \mathbb{R}$ at $w \in \mathbb{R}$ is defined as the derivative of f' at w , that is,

$$f^{(2)}(w) := \lim_{\delta \rightarrow 0} \frac{f'(w + \delta) - f'(w)}{\delta},$$

provided that the limit is well-defined. If the second derivative of a function f is well-defined at w , the function is said to be **twice differentiable** at w . The second derivative is also denoted f'' .

If f has a small second derivative at a given w , the derivative around w is almost constant. That is, the function behaves like a line around w , as illustrated in Fig. 2.7. Hence, the second derivative is usually interpreted as the **curvature** of the function at a given point.

2.4.2 Second directional derivatives

For a multi-input function $f : \mathbb{R}^P \rightarrow \mathbb{R}$, we saw that the directional derivative encodes infinitesimal variations of f along a given direction. To analyze the second derivative, the curvature of the function at a given point \mathbf{w} , we can consider the variations along a pair of directions, as defined below.

Definition 2.15 (Second directional derivative). The **second directional derivative** of $f : \mathbb{R}^P \rightarrow \mathbb{R}$ at $\mathbf{w} \in \mathbb{R}^P$ along $\mathbf{v}, \mathbf{v}' \in \mathbb{R}^P$ is defined as the directional derivative of $\mathbf{w} \mapsto \partial f(\mathbf{w})[\mathbf{v}]$ along \mathbf{v}' , that is,

$$\partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{v}'] := \lim_{\delta \rightarrow 0} \frac{\partial f(\mathbf{w} + \delta \mathbf{v}')[\mathbf{v}] - \partial f(\mathbf{w})[\mathbf{v}]}{\delta},$$

provided that $\partial f(\mathbf{w})[\mathbf{v}]$ is well-defined around \mathbf{w} and that the limit exists.

Of particular interest are the variations of a function around the canonical directions: the **second partial derivatives**, defined as

$$\partial_{ij}^2 f(\mathbf{w}) := \partial^2 f(\mathbf{w})[\mathbf{e}_i, \mathbf{e}_j]$$

for $\mathbf{e}_i, \mathbf{e}_j$ the i^{th} and j^{th} canonical directions in \mathbb{R}^P , respectively. In Leibniz notation, the second partial derivatives are denoted

$$\partial_{ij}^2 f(\mathbf{w}) = \frac{\partial^2 f(\mathbf{w})}{\partial w_i \partial w_j}.$$

2.4.3 Hessians

For a multi-input function, twice differentiability is simply defined as the differentiability of any directional derivative $\partial f(\mathbf{w})[\mathbf{v}]$ w.r.t. \mathbf{w} .

Definition 2.16 (Twice differentiability). A function $f : \mathbb{R}^P \rightarrow \mathbb{R}$ is twice differentiable at $\mathbf{w} \in \mathbb{R}^P$ if it is differentiable and $\partial f : \mathbb{R}^P \rightarrow (\mathbb{R}^P \rightarrow \mathbb{R})$ is also differentiable at \mathbf{w} .

As a result, the second directional derivative is a bilinear form.

Definition 2.17 (Bilinear map, bilinear form). A function $b : \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{R}^M$ is a **bilinear map** if $b[\mathbf{v}, \cdot] : \mathbb{R}^P \rightarrow \mathbb{R}^M$ is linear for any \mathbf{v} and $b[\cdot, \mathbf{v}']$ is linear for any \mathbf{v}' . That is,

$$b[\mathbf{v}, \mathbf{v}'] = \sum_{i=1}^P v_i b[\mathbf{e}_i, \mathbf{v}'] = \sum_{i=1}^P \sum_{j=1}^P v_i v'_j b[\mathbf{e}_i, \mathbf{e}_j],$$

for $\mathbf{v} = \sum_{i=1}^P v_i \mathbf{e}_i$ and $\mathbf{v}' = \sum_{i=1}^P v'_i \mathbf{e}_i$. A bilinear map with values in \mathbb{R} , $b : \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{R}$, is called a **bilinear form**.

The second partial derivatives are gathered in the **Hessian** and the second directional derivatives can be computed from it.

Definition 2.18 (Hessian). The **Hessian** of a twice differentiable function $f : \mathbb{R}^P \rightarrow \mathbb{R}$ at \mathbf{w} is the $P \times P$ matrix gathering all second partial derivatives,

$$\nabla^2 f(\mathbf{w}) := \begin{pmatrix} \partial_{11}f(\mathbf{w}) & \dots & \partial_{1P}f(\mathbf{w}) \\ \vdots & \ddots & \vdots \\ \partial_{P1}f(\mathbf{w}) & \dots & \partial_{PP}f(\mathbf{w}) \end{pmatrix} \in \mathbb{R}^{P \times P},$$

provided that all second partial derivatives are well-defined.

The second directional derivative at \mathbf{w} is bilinear in any directions $\mathbf{v} = \sum_{i=1}^P v_i \mathbf{e}_i$ and $\mathbf{v}' = \sum_{i=1}^P v'_i \mathbf{e}_i$. Therefore,

$$\partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{v}'] = \sum_{i,j=1}^P v_i v'_j \partial^2 f(\mathbf{w})[\mathbf{e}_i, \mathbf{e}_j] = \langle \mathbf{v}, \nabla^2 f(\mathbf{w}) \mathbf{v}' \rangle.$$

Given the gradient of f , the Hessian is equivalent to the transpose of the Jacobian of the gradient. By slightly generalizing the notation ∇ to denote the transpose of the Jacobian of a function (which matches its definition for single-output functions), we have that the Hessian can be expressed as $\nabla^2 f(\mathbf{w}) = \nabla(\nabla f)(\mathbf{w})$, which justifies its notation.

Similarly as for the differentiability of a function f , twice differentiability of f at \mathbf{w} is equivalent to having the second partial derivatives not only defined but also continuous in a neighborhood of \mathbf{w} . Remarkably, by requiring twice differentiability, i.e., continuous second partial derivatives, the Hessian is guaranteed to be symmetric (Schwarz, 1873).

Proposition 2.10 (Symmetry of the Hessian). If a function $f : \mathbb{R}^P \rightarrow \mathbb{R}$ is twice differentiable at \mathbf{w} , then its Hessian $\nabla^2 f(\mathbf{w})$ is symmetric, that is, $\partial_{ij}^2 f(\mathbf{w}) = \partial_{ji}^2 f(\mathbf{w})$ for any $i, j \in \{1, \dots, P\}$.

The symmetry of the Hessian means that it can alternatively be written as $\nabla^2 f(\mathbf{w}) = (\partial_{ji}^2 f(\mathbf{w}))_{i,j=1}^P = \partial(\nabla f)(\mathbf{w})$, i.e., the Jacobian of the gradient of f .

2.4.4 Hessian-vector products

Similarly to the Jacobian, we can exploit the formal definition of the Hessian as a bilinear form to extend its definition to Euclidean spaces. In particular, we can define the notion of Hessian-vector product.

Definition 2.19 (Hessian-vector product). If a function $f : \mathcal{E} \rightarrow \mathbb{R}$ defined on a Euclidean space \mathcal{E} with inner product $\langle \cdot, \cdot \rangle$, is twice differentiable at $\mathbf{w} \in \mathcal{E}$, then for any $\mathbf{v} \in \mathcal{E}$, there exists $\mathbf{v} \mapsto \nabla^2 f(\mathbf{w})[\mathbf{v}]$, called the **Hessian-vector product** (HVP) of f at \mathbf{w} along \mathbf{v} , such that for any $\mathbf{v}' \in \mathcal{E}$,

$$\partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{v}'] = \langle \mathbf{v}', \nabla^2 f(\mathbf{w})[\mathbf{v}] \rangle.$$

In particular for $\mathcal{E} = \mathbb{R}^P$, the HVP is $\nabla^2 f(\mathbf{w})[\mathbf{v}] = (\partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{e}_i])_{i=1}^P$.

From an autodiff point of view, the HVP can be implemented in four different ways, as explained in Section 9.1.

2.4.5 Second-order Jacobians

The previous definitions naturally extend to multi-output functions $f : \mathcal{E} \rightarrow \mathcal{F}$, where $f := (f_1, \dots, f_M)$, $f_j : \mathcal{E} \rightarrow \mathcal{F}_j$ and $\mathcal{F} := \mathcal{F}_1 \times \dots \times \mathcal{F}_M$. The second directional derivative is defined by gathering the second derivatives of each coordinate's function. That is, for $\mathbf{w}, \mathbf{v}, \mathbf{v}' \in \mathcal{E}$,

$$\partial f(\mathbf{w})[\mathbf{v}, \mathbf{v}'] = (\partial f_j(\mathbf{w})[\mathbf{v}, \mathbf{v}'])_{j=1}^M \in \mathcal{F}.$$

The function f is twice differentiable if and only if all its coordinates are twice differentiable. The second directional derivative is then a **bilinear**

map. We can then compute second directional derivatives as

$$\partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{v}'] = \sum_{i,j=1}^P v_i v'_j \partial^2 f(\mathbf{w})[\mathbf{e}_i, \mathbf{e}_j] = (\langle \mathbf{v}, \nabla^2 f_j(\mathbf{w}) \mathbf{v}' \rangle)_{j=1}^M.$$

When $\mathcal{E} = \mathbb{R}^P$ and $\mathcal{F}_j = \mathbb{R}$, so that $\mathcal{F} = \mathbb{R}^M$, the bilinear map can be materialized as a tensor

$$\partial^2 f(\mathbf{w}) = (\partial^2 f(\mathbf{w})[\mathbf{e}_i, \mathbf{e}_j])_{i,j=1}^P \in \mathbb{R}^{M \times P \times P},$$

the “second-order Jacobian” of f . However, similarly to the Hessian, it is usually more convenient to apply the bilinear map to prescribed vectors \mathbf{v} and \mathbf{v}' than to materialize the second partial derivatives as a tensor.

2.5 Higher-order differentiation

2.5.1 Higher-order derivatives

Derivatives can be extended to any order. Formally, the n^{th} derivative can be defined inductively as follows for a single-input, single-output function.

Definition 2.20 (n^{th} order derivative). The n^{th} derivative $f^{(n)}$ of a function $f : \mathbb{R} \rightarrow \mathbb{R}$ at $w \in \mathbb{R}$ is defined as

$$f^{(n)}(w) := (f^{(n-1)})'(w) = \lim_{\delta \rightarrow 0} \frac{f^{(n-1)}(w + \delta) - f^{(n-1)}(w)}{\delta}$$

provided that $f^{(n-1)}$ is differentiable around w and that the limit exists. In such a case, the function is said to be n times differentiable at w .

2.5.2 Higher-order directional derivatives

For a multi-input function f , we can naturally extend the notion of directional derivative as follows.

Definition 2.21 (n^{th} order directional derivative). The n^{th} directional derivative of $f : \mathbb{R}^P \rightarrow \mathbb{R}$ at $\mathbf{w} \in \mathbb{R}^P$ along $\mathbf{v}_1, \dots, \mathbf{v}_n$ is defined as

$$\begin{aligned} \partial^n f(\mathbf{w})[\mathbf{v}_1, \dots, \mathbf{v}_n] &= \partial(\partial^{n-1} f(\mathbf{w})[\mathbf{v}_1, \dots, \mathbf{v}_{n-1}])[\mathbf{v}_n] \\ &= \lim_{\delta \rightarrow 0} \frac{\partial f(\mathbf{w} + \delta \mathbf{v}_n)[\mathbf{v}_1, \dots, \mathbf{v}_{n-1}] - \partial f(\mathbf{w})[\mathbf{v}_1, \dots, \mathbf{v}_{n-1}]}{\delta} \end{aligned}$$

A multi-input function f is n -times differentiable if it is $n - 1$ differentiable and its $n - 1$ directional derivative along any direction is differentiable. As a consequence the n^{th} directional derivative is a **multilinear form**.

Definition 2.22 (Multilinear map, multilinear form). A function $c : \otimes_{i=1}^n \mathbb{R}^P \rightarrow \mathbb{R}^M$ is a **multilinear map** if it is linear in each coordinate given all others fixed, that is, if $\mathbf{v}_j \mapsto c[\mathbf{v}_1, \dots, \mathbf{v}_j, \dots, \mathbf{v}_n]$ is linear in \mathbf{v}_j for any $j \in [n]$. It is a **multilinear form** if it has values in \mathbb{R} .

The n^{th} order directional derivative is then given by

$$\partial^n f(\mathbf{w})[\mathbf{v}_1, \dots, \mathbf{v}_n] = \sum_{i_1, \dots, i_n=1}^P v_{1,i_1} \dots v_{n,i_n} \partial^n f(\mathbf{w})[\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_n}].$$

The n^{th} order partial derivatives can be materialized as an n^{th} order tensor

$$\nabla^n f(\mathbf{w}) = (\partial^n f(\mathbf{w})[\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_n}])_{i_1, \dots, i_n=1}^P \in \mathbb{R}^{P \times \dots \times P}.$$

2.5.3 Higher-order Jacobians

All above definitions extend directly to the case of multi-output functions $f : \mathcal{E} \rightarrow \mathcal{F}$, where $\mathcal{F} := \mathcal{F}_1 \times \dots \times \mathcal{F}_M$. The n^{th} directional derivatives are then

$$\partial^n f(\mathbf{w})[\mathbf{v}_1, \dots, \mathbf{v}_n] = (\partial^n f_j(\mathbf{w})[\mathbf{v}_1, \dots, \mathbf{v}_n])_{j=1}^M.$$

The function f is then n times differentiable if it is $n - 1$ differentiable and its $n - 1$ directional derivative along any direction is differentiable.

As a consequence, the n^{th} directional derivative is a **multilinear map**. The n^{th} directional derivative can be decomposed into partial derivatives as

$$\partial^n f(\mathbf{w})[\mathbf{v}_1, \dots, \mathbf{v}_n] = \sum_{i_1, \dots, i_n=1}^P v_{1,i_1} \dots v_{n,i_n} \partial^n f(\mathbf{w})[\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_n}].$$

When $\mathcal{E} = \mathbb{R}^P$ and $\mathcal{F} = \mathbb{R}^M$, the n^{th} order partial derivatives can be materialized by an $n + 1^{\text{th}}$ order tensor

$$\partial^n f(\mathbf{w}) = (\partial^n f_j(\mathbf{w})[\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_n}])_{j=1, i_1, \dots, i_n=1}^{M, P, \dots, P} \in \mathbb{R}^{M \times P \times \dots \times P}.$$

2.5.4 Taylor expansions

With Landau's little o notation, we have seen that if a function is differentiable, it is approximated by a linear function in \mathbf{v} ,

$$f(\mathbf{w} + \mathbf{v}) = f(\mathbf{w}) + \langle \nabla f(\mathbf{w}), \mathbf{v} \rangle + o(\|\mathbf{v}\|_2).$$

Such an expansion of the function up to its first derivative is called the **first-order Taylor expansion** of f around \mathbf{w} .

If the function f is twice differentiable, we can approximate it by a quadratic in \mathbf{v} , leading to the **second-order Taylor expansion** of f around \mathbf{w} ,

$$f(\mathbf{w} + \mathbf{v}) = f(\mathbf{w}) + \langle \nabla f(\mathbf{w}), \mathbf{v} \rangle + \frac{1}{2} \langle \mathbf{v}, \nabla^2 f(\mathbf{w}) \mathbf{v} \rangle + o(\|\mathbf{v}\|_2^2).$$

Compared to the first-order Taylor approximation, it is naturally more accurate around \mathbf{w} , as reflected by the fact that $\|\mathbf{v}\|_2^3 \leq \|\mathbf{v}\|_2^2$ for $\|\mathbf{v}\|_2 \leq 1$.

More generally, we can build the **n^{th} order Taylor expansion** of a n times differentiable function $f: \mathbb{R}^P \rightarrow \mathbb{R}^M$ around $\mathbf{w} \in \mathbb{R}^P$ by

$$\begin{aligned} f(\mathbf{w} + \mathbf{v}) &= f(\mathbf{w}) + \partial f(\mathbf{w})[\mathbf{v}] + \frac{1}{2} \partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{v}] + \dots \\ &\quad + \frac{1}{n!} \partial^n f(\mathbf{w})[\underbrace{\mathbf{v}, \dots, \mathbf{v}}_{n \text{ times}}] + o(\|\mathbf{v}\|_2^n). \end{aligned}$$

Note that, using the change of variable $\mathbf{w}' = \mathbf{w} + \mathbf{v} \iff \mathbf{v} = \mathbf{w}' - \mathbf{w}$, it is often convenient to write the n^{th} Taylor expansion of $f(\mathbf{w}')$ around

\mathbf{w} as

$$f(\mathbf{w}') = f(\mathbf{w}) + \sum_{j=1}^n \frac{1}{j!} \partial^j f(\mathbf{w}) [\underbrace{\mathbf{w}' - \mathbf{w}, \dots, \mathbf{w}' - \mathbf{w}}_{j \text{ times}}] + o(\|\mathbf{w}' - \mathbf{w}\|_2^n).$$

Taylor expansions will prove useful in Chapter 7 for computing derivatives by finite differences.

2.6 Differential geometry

In this chapter, we progressively generalized the notion of derivative from real numbers to vectors and variables living in a linear space (a.k.a. vector space), either finite dimensional or infinite dimensional. We can further generalize these notions by considering a local notion of linearity. This is formalized by smooth manifolds in **differential geometry**, whose terminology is commonly adopted in the automatic differentiation literature and software. In this section, we give a brief overview of derivatives on smooth manifolds (simply referred to as manifolds), and refer to Boumal (2023) for a complete introduction.

2.6.1 Differentiability on manifolds

Essentially, a manifold is a set that can be locally approximated by a Euclidean space. The most common example is a sphere like the Earth. Seen from the Moon, the Earth is not a plane, but locally, at a human level, it can be seen as a flat surface. Euclidean spaces are also trivial examples of manifolds. A formal characterization of the sphere as a manifold is presented in Example 2.9. For now, we may think of a “manifold” as some set (e.g., the sphere) contained in some ambient Euclidean space; note however that manifolds can be defined generally without being contained in a Euclidean space (Boumal, 2023, Chapter 8). Differentiability in manifolds is simply inherited from the notion of differentiability in the ambient Euclidean space.

Definition 2.23 (Differentiability of restricted functions). Let \mathcal{M} and \mathcal{N} be manifolds. A function $f : \mathcal{M} \rightarrow \mathcal{N}$ defined from $\mathcal{M} \subseteq \mathcal{E}$ to $\mathcal{N} \subseteq \mathcal{F}$, with \mathcal{E} and \mathcal{F} Euclidean spaces, is differentiable if f is

the restriction of a differentiable function $\bar{f} : \mathcal{E} \rightarrow \mathcal{F}$, so that f coincides with \bar{f} on \mathcal{M} .

Our objective is to formalize the directional derivatives and gradients for functions defined on manifolds. This formalization leads to the definitions of tangent spaces and cotangent spaces, and the associated generalizations of JVP and VJP operators as pushforward and pullback operators, respectively.

2.6.2 Tangent spaces and pushforward operators

To generalize the notion of directional derivatives of a function f , the one property we want to preserve is the chain rule. Rather than starting from the variations of f at a given point along a direction, we start with the variations of f along curves. Namely, on a manifold like the sphere \mathcal{S}^P in \mathbb{R}^P , we can look at curves $\alpha : \mathbb{R} \rightarrow \mathcal{S}^P$ passing through $\mathbf{w} \in \mathcal{S}^P$ at time 0, that is, $\alpha(0) = \mathbf{w}$. For single-input functions like α , we denote for simplicity $\alpha'(0) := (\alpha'_1(0), \dots, \alpha'_P(0))$. The directional derivative of a function f must typically serve to define the derivative of $f \circ \alpha$ at 0, such that $(f \circ \alpha)'(0) = \partial f(\mathbf{w})[\alpha'(0)]$. In the case of the sphere, as illustrated in Fig. 2.8, the derivative $\alpha'(0)$ of a curve α passing through a point \mathbf{w} is always **tangent** to the sphere at \mathbf{w} . The tangent plane to the sphere at \mathbf{w} then captures all possible relevant vectors to pass to the JVP we are building. To define the directional derivative of a function f on a manifold, we therefore restrict ourselves to an operator defined on the **tangent space** $\mathcal{T}_{\mathbf{w}}\mathcal{M}$, whose definition below is simplified for our purposes.

Definition 2.24 (Tangent space). The **tangent space** of a manifold \mathcal{M} at $\mathbf{w} \in \mathcal{M}$ is defined as

$$\mathcal{T}_{\mathbf{w}}\mathcal{M} := \{\mathbf{v} = \alpha'(0) \text{ for any } \alpha : \mathbb{R} \rightarrow \mathcal{M} \text{ differentiable s.t. } \alpha(0) = \mathbf{w}\}.$$

In the case of the sphere in Fig. 2.8, the tangent space is a plane, that is, a Euclidean space. This property is generally true: tangent spaces are Euclidean spaces, enabling us to define directional derivatives as linear maps. Now, if f is differentiable and goes from a manifold \mathcal{M} to a manifold \mathcal{N} , then $f \circ \alpha$ is a differentiable curve in \mathcal{N} . Therefore,

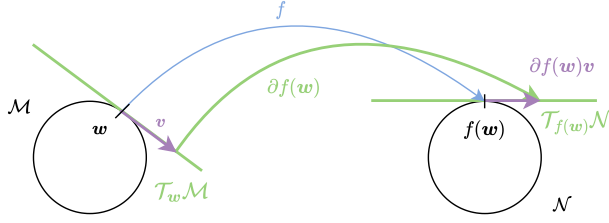


Figure 2.8: A differentiable function f defined from a sphere \mathcal{M} to a sphere \mathcal{N} defines a push-forward operator that maps tangent vectors (derivatives of functions on the sphere passing through w) in the tangent space $\mathcal{T}_w \mathcal{M}$ to tangent vectors of \mathcal{N} at $f(w)$ in the tangent space $\mathcal{T}_{f(w)} \mathcal{N}$.

$(f \circ \alpha)'(0)$ is the derivative of a curve passing through $f(w)$ at 0 and is tangent to \mathcal{N} at $f(w)$. Hence, the directional derivative of $f : \mathcal{M} \rightarrow \mathcal{N}$ at w can be defined as a function from the tangent space $\mathcal{T}_w \mathcal{M}$ of \mathcal{M} at w onto the tangent space $\mathcal{T}_{f(w)} \mathcal{N}$ of \mathcal{N} at $f(w)$. Overall, we built the directional derivative (JVP) by considering how a composition of f with any curve α pushes forward the derivative of α into the derivative of $f \circ \alpha$. The resulting JVP is called a **pushforward operator** in differentiable geometry.

Definition 2.25 (Pushforward operator). Given two manifolds \mathcal{M} and \mathcal{N} , the **pushforward operator** of a differentiable function $f : \mathcal{M} \rightarrow \mathcal{N}$ at $w \in \mathcal{M}$ is the linear map $\partial f(w) : \mathcal{T}_w \mathcal{M} \rightarrow \mathcal{T}_{f(w)} \mathcal{N}$ defined by

$$\partial f(w)[v] := (f \circ \alpha)'(0),$$

for any $v \in \mathcal{T}_w \mathcal{M}$ such that $v = \alpha'(0)$, where $\alpha : \mathbb{R} \rightarrow \mathcal{M}$ is a differentiable curve passing through w at 0, i.e., $\alpha(0) = w$.

2.6.3 Cotangent spaces and pullback operators

To generalize the JVP, we composed $f : \mathcal{M} \rightarrow \mathcal{N}$ with any single-input function $\alpha : \mathbb{R} \rightarrow \mathcal{M}$ giving values on the manifold. The derivative of any such α is then pushed forward from $\mathcal{T}_w \mathcal{M}$ to $\mathcal{T}_{f(w)} \mathcal{N}$ by the action of f . To define the VJP, we take a symmetric approach. We consider all single-output differentiable functions $\beta : \mathcal{N} \rightarrow \mathbb{R}$ defined on $y \in \mathcal{N}$ with

$\mathbf{y} = f(\mathbf{w})$ for some $\mathbf{w} \in \mathcal{M}$. We then want to pull back the derivatives of β when precomposing it by f . Therefore, the space on which the VJP acts is the space of directional derivatives of any $\beta : \mathcal{N} \rightarrow \mathbb{R}$ at \mathbf{y} , defining the **cotangent space**.

Definition 2.26 (Cotangent space). The **cotangent space** of a manifold \mathcal{N} at $\mathbf{y} \in \mathcal{N}$ is defined as

$$\begin{aligned}\mathcal{T}_{\mathbf{y}}^* \mathcal{N} &= \{u = \partial\beta(\mathbf{y}) \text{ for any } \beta : \mathcal{N} \rightarrow \mathbb{R} \text{ differentiable}\} \\ &= \{u : \mathcal{T}_{\mathbf{y}} \mathcal{N} \rightarrow \mathbb{R} \text{ for any linear map } u\},\end{aligned}$$

Note that elements of the cotangent space are linear mappings, not vectors. This distinction is important to define the pullback operator as an operator on functions as done in measure theory. From a linear algebra viewpoint, the cotangent space is exactly the **dual space** of $\mathcal{T}_{\mathbf{y}} \mathcal{N}$, that is, the set of linear maps from $\mathcal{T}_{\mathbf{y}} \mathcal{N}$ to \mathbb{R} , called **linear forms**. As $\mathcal{T}_{\mathbf{y}} \mathcal{N}$ is a Euclidean space, its dual space $\mathcal{T}_{\mathbf{y}}^* \mathcal{N}$ is also a Euclidean space. The **pullback** operator is then defined as the operator that gives access to directional derivatives of $\beta \circ f$ given the directional derivative of β at $f(\mathbf{w})$.

Definition 2.27 (Pullback operator). Given two manifolds \mathcal{M} and \mathcal{N} , the **pullback operator** of a differentiable function $f : \mathcal{M} \rightarrow \mathcal{N}$ at $\mathbf{w} \in \mathcal{M}$ is the linear map $\partial f(\mathbf{w})^* : \mathcal{T}_{f(\mathbf{w})}^* \mathcal{N} \rightarrow \mathcal{T}_{\mathbf{w}}^* \mathcal{M}$ defined by

$$\partial f(\mathbf{w})^* u := \partial(\beta \circ f)(\mathbf{w}),$$

for any $u \in \mathcal{T}_{f(\mathbf{w})}^* \mathcal{N}$ such that $\partial\beta(f(\mathbf{w})) = u$, for a differentiable function $\beta : \mathcal{N} \rightarrow \mathbb{R}$.

Contrary to the pushforward operator that acts on vectors, the pullback operator acts on linear forms. Hence, the slight difference in notation between $\partial f(\mathbf{w})^*$ and $\partial f(\mathbf{w})^*$, the adjoint operator of $\partial f(\mathbf{w})$. To properly define the adjoint operator $\partial f(\mathbf{w})^*$, we need a notion of inner product. Since tangent spaces are Euclidean spaces, we can define an inner product $\langle \cdot, \cdot \rangle_{\mathbf{w}}$ for each $\mathcal{T}_{\mathbf{w}} \mathcal{M}$ and $\mathbf{w} \in \mathcal{M}$, making \mathcal{M} a **Riemannian manifold**. Equipped with these inner products, the cotangent space can be identified with the tangent space, and we can define gradients.

Function	f	$\mathcal{M} \rightarrow \mathcal{N}$
Push-forward	$\partial f(\mathbf{w})$	$\mathcal{T}_{\mathbf{w}}\mathcal{M} \rightarrow \mathcal{T}_{f(\mathbf{w})}\mathcal{N}$
Pullback	$\partial f(\mathbf{w})^*$	$\mathcal{T}_{f(\mathbf{w})}^*\mathcal{N} \rightarrow \mathcal{T}_{\mathbf{w}}^*\mathcal{M}$
Adjoint of pushforward	$\partial f(\mathbf{w})^*$	$\mathcal{T}_{f(\mathbf{w})}\mathcal{N} \rightarrow \mathcal{T}_{\mathbf{w}}\mathcal{M}$

Table 2.1: For a differentiable function f defined from a manifold \mathcal{M} onto a manifold \mathcal{N} , the JVP is generalized with the notion of pushforward $\partial f(\mathbf{w})$. The counterpart of the pushforward is the pullback operation $\partial f(\mathbf{w})^*$ that acts on linear forms in the tangent spaces. For Riemannian manifolds, the pullback operation can be identified with the adjoint operator $\partial f(\mathbf{w})^*$ of the pushforward operator as any linear form is represented by a vector.

Definition 2.28 (Gradients in Riemannian manifolds). Let \mathcal{M} be a Riemannian manifold equipped with inner products $\langle \cdot, \cdot \rangle_{\mathbf{w}}$. For any **cotangent vector** $u \in \mathcal{T}_{\mathbf{w}}^*\mathcal{M}$, with $\mathbf{w} \in \mathcal{M}$, there exists a unique **tangent vector** $\mathbf{u} \in \mathcal{T}_{\mathbf{w}}\mathcal{M}$ such that

$$\forall \mathbf{v} \in \mathcal{T}_{\mathbf{w}}\mathcal{M}, \quad u[\mathbf{v}] = \langle \mathbf{u}, \mathbf{v} \rangle_{\mathbf{w}}.$$

In particular for any differentiable function $f : \mathcal{M} \rightarrow \mathbb{R}$, we can define the **gradient** of f as the unique tangent vector $\nabla f(\mathbf{w}) \in \mathcal{T}_{\mathbf{w}}\mathcal{M}$ such that

$$\forall \mathbf{v} \in \mathcal{T}_{\mathbf{w}}\mathcal{M}, \quad \partial f(\mathbf{w})[\mathbf{v}] = \langle \nabla f(\mathbf{w}), \mathbf{v} \rangle.$$

Therefore, rather than pulling back directional derivatives, we can pull back gradients. The corresponding operator is then naturally the adjoint $\partial f(\mathbf{w})^*$ of the pushforward operator. Namely, given two Riemannian manifolds \mathcal{M} and \mathcal{N} , and a differentiable function $f : \mathcal{M} \rightarrow \mathcal{N}$, we have

$$(\partial f(\mathbf{w})^*u)[\mathbf{v}] = \langle \partial f(\mathbf{w})^*[\mathbf{u}], \mathbf{v} \rangle \text{ for any } \mathbf{v} \in \mathcal{T}_{\mathbf{w}}\mathcal{M}$$

for $u = \langle \cdot, \mathbf{u} \rangle \in \mathcal{T}_{f(\mathbf{w})}^*\mathcal{N}$ represented by $\mathbf{u} \in \mathcal{T}_{f(\mathbf{w})}\mathcal{N}$.

Example 2.9 (The sphere as a manifold). The sphere \mathcal{S}^P in \mathbb{R}^P is defined as the set of points $\mathbf{w} \in \mathbb{R}^P$, satisfying $c(\mathbf{w}) := \langle \mathbf{w}, \mathbf{w} \rangle - 1 = 0$, with JVP $\partial c(\mathbf{w})[\mathbf{v}] = 2\langle \mathbf{w}, \mathbf{v} \rangle$.

For any $\mathbf{v} = (v_1, \dots, v_{P-1}) \in \mathbb{R}^{P-1}$ close enough to a point \mathbf{w} on the sphere, we can define $\psi_1(\mathbf{v}) := \sqrt{1 - \langle \mathbf{v}, \mathbf{v} \rangle}$ such that $\boldsymbol{\psi}(\mathbf{v}) := (v_1, \dots, v_{P-1}, \psi_1(\mathbf{v}))$ satisfies $\langle \boldsymbol{\psi}(\mathbf{v}), \boldsymbol{\psi}(\mathbf{v}) \rangle = 1$, that is $c(\boldsymbol{\psi}(\mathbf{v})) = 1$. With the help of the mapping $\boldsymbol{\psi}^{-1}$ from a neighborhood of \mathbf{w} in the sphere to \mathbb{R}^{P-1} , we can locally see the sphere as a space of dimension $P - 1$.

The tangent space can be naturally characterized in terms of the constraining function c . Namely, the curve $\alpha : \mathbb{R} \rightarrow \mathcal{S}$ such that $\alpha(0) = \mathbf{w}$ satisfies for any $\delta \in \mathbb{R}$, $c(\alpha(\delta)) = \mathbf{0}$. Hence, differentiating the implicit equation, we have

$$(c \circ \alpha)'(0) = \partial c(\mathbf{w})[\alpha'(0)].$$

That is, $\alpha'(0)$ is in the null space of $\partial c(\mathbf{w})$, denoted

$$\text{Null}(\partial c(\mathbf{w})) := \{\mathbf{v} \in \mathbb{R}^P : \partial c(\mathbf{w})[\mathbf{v}] = 0\}.$$

The tangent space of \mathcal{S} at \mathbf{w} is then

$$\begin{aligned} \mathcal{T}_{\mathbf{w}}\mathcal{M} &= \text{Null}(2\langle \mathbf{w}, \cdot \rangle) \\ &= \{\mathbf{v} \in \mathbb{R}^P : \langle \mathbf{w}, \mathbf{v} \rangle = 0\} \end{aligned}$$

We naturally recover that the tangent space is a Euclidean space of dimension $P - 1$, defined as the set of points orthogonal to \mathbf{w} .

2.7 Generalized derivatives

While we largely focus on differentiable functions in this book, it is important to characterize non-differentiable functions. We distinguish here two cases: continuous functions and non-continuous functions. For the former case, there exist generalizations of the notion of directional derivative, gradient and Jacobian, presented below. For non-continuous functions, even if derivatives exist almost everywhere, they may be uninformative. For example, piecewise-constant functions, encountered in e.g. control flows (Chapter 5), are almost everywhere differentiable but with zero derivatives. In such cases, surrogate functions can be defined to ensure the differentiability of a program (Part IV).

2.7.1 Rademacher's theorem

We first recall the definition of (locally) Lipschitz continuous function.

Definition 2.29 ((Locally) Lipschitz continuous function). A function $f : \mathcal{E} \rightarrow \mathcal{F}$, is Lipschitz continuous if there exists $C \geq 0$ such that for any $\mathbf{x}, \mathbf{y} \in \mathcal{E}$,

$$\|f(\mathbf{x}) - f(\mathbf{y})\| \leq C\|\mathbf{x} - \mathbf{y}\|.$$

A function $f : \mathcal{E} \rightarrow \mathcal{F}$ is locally Lipschitz continuous if for any $\mathbf{x} \in \mathcal{E}$, there exists a neighborhood \mathcal{U} of \mathbf{x} such that f restricted to \mathcal{U} is Lipschitz continuous.

Rademacher's theorem (Rademacher, 1919) then ensures that f is differentiable almost everywhere.

Proposition 2.11 (Rademacher's theorem). Let \mathcal{E} and \mathcal{F} denote Euclidean spaces. If $f : \mathcal{E} \rightarrow \mathcal{F}$ is locally Lipschitz-continuous, then f is almost everywhere differentiable, that is, the set of points in \mathcal{E} at which f is not differentiable is of (Lebesgue) measure zero.

See also Morrey Jr (2009) for a standard proof.

2.7.2 Clarke derivatives

Rademacher's theorem hints that the definitions of directional derivatives, gradients and Jacobians may be generalized to locally Lipschitz continuous functions. This is what Clarke (1975) did in his seminal work, which laid the foundation of **nonsmooth analysis**. The first building block is a notion of generalized directional derivative.

Definition 2.30 (Clarke generalized directional derivative). The **Clarke generalized directional derivative** of a locally Lipschitz continuous function $f : \mathcal{E} \rightarrow \mathbb{R}$ at $\mathbf{w} \in \mathcal{E}$ in the direction $\mathbf{v} \in \mathcal{E}$ is

$$\partial_C f(\mathbf{w})[\mathbf{v}] := \limsup_{\substack{\mathbf{u} \rightarrow \mathbf{w} \\ \delta \searrow 0}} \frac{f(\mathbf{u} + \delta \mathbf{v}) - f(\mathbf{u})}{\delta},$$

provided that the limit exists, where $\delta \searrow 0$ means that δ approaches

0 by non-negative values, and where the limit superior is defined as

$$\limsup_{\mathbf{x} \rightarrow \mathbf{a}} f(\mathbf{x}) := \lim_{\varepsilon \rightarrow 0} \sup \{f(\mathbf{x}) : \mathbf{x} \in B(\mathbf{a}, \varepsilon) \setminus \{\mathbf{a}\}\}$$

for $B(\mathbf{a}, \varepsilon) := \{\mathbf{x} \in \mathcal{E} : \|\mathbf{x} - \mathbf{a}\| \leq \varepsilon\}$ the ball centered at \mathbf{a} of radius ε .

There are two differences with the usual definition of a directional derivative: (i) we consider slopes of the function in a neighborhood of the point rather than at the given point, (ii) we take a limit superior rather than a usual limit. The first point is rather natural in the light of Rademacher's theorem: we can properly characterize variations on points where the function is differentiable, therefore we may take the limits of these slopes as a candidate slope for the point of interest. The second point is more technical but essential: it allows us to characterize the directional derivative as the supremum of some linear forms (Clarke *et al.*, 2008). These linear forms in turn define a set of generalized gradients (Clarke *et al.*, 2008, Chapter 2).

Definition 2.31 (Clarke generalized gradient). A **Clarke generalized gradient** of a locally Lipschitz function $f : \mathcal{E} \rightarrow \mathbb{R}$ at $\mathbf{w} \in \mathcal{E}$ is a point $\mathbf{g} \in \mathcal{E}$ such that $\forall \mathbf{v} \in \mathcal{E}$

$$\partial f(\mathbf{w})[\mathbf{v}] \geq \langle \mathbf{g}, \mathbf{v} \rangle.$$

The set of Clarke generalized gradients is called the **Clarke subdifferential** of f at \mathbf{w} .

Definition 2.30 and Definition 2.31 can be used in non-Euclidean spaces, such as Banach or Hilbert spaces (Clarke *et al.*, 2008). In Euclidean spaces, the Clarke generalized gradients can be characterized more simply thanks to Rademacher's theorem (Clarke *et al.*, 2008, Theorem 8.1). Namely, as shown below, they can be defined as a convex combination of limits of gradients of f evaluated at a sequence in $\mathcal{E} \setminus \Omega$ that converges to \mathbf{w} .

Proposition 2.12 (Characterization of Clarke generalized gradients).

Let $f : \mathcal{E} \rightarrow \mathbb{R}$ be a locally Lipschitz continuous and denote Ω the set of points at which f is not differentiable (Proposition 2.11). An element $\mathbf{g} \in \mathcal{E}$ is a Clarke generalized gradient of f at $\mathbf{w} \in \mathcal{E}$ if and only if

$$\mathbf{g} \in \text{conv} \left(\left\{ \lim_{n \rightarrow +\infty} \nabla f(\mathbf{v}_n) : (\mathbf{v}_n)_{n=1}^{+\infty} \text{ s.t. } \mathbf{v}_n \in \mathcal{E} \setminus \Omega, \mathbf{v}_n \xrightarrow{n \rightarrow +\infty} \mathbf{w} \right\} \right).$$

In the above, the convex hull of a set $S \subseteq \mathcal{E}$, the set of convex combinations of elements of S , is denoted

$$\text{conv}(S) := \{ \lambda_1 \mathbf{s}_1 + \dots + \lambda_m \mathbf{s}_m : m \in \mathbb{N}, \lambda_i \geq 0, \sum_{i=1}^m \lambda_i = 1, \mathbf{s}_i \in S \}.$$

The Jacobian of a function $f : \mathcal{E} \rightarrow \mathcal{F}$ between two Euclidean spaces can be generalized similarly (Clarke *et al.*, 2008, Section 3.3).

Definition 2.32 (Clarke generalized Jacobian). Let $f : \mathcal{E} \rightarrow \mathcal{F}$ be a locally Lipschitz continuous and denote Ω the set of points at which f is not differentiable (Proposition 2.11). A **Clarke generalized Jacobian** of f at $\mathbf{w} \in \mathcal{E}$ is an element \mathbf{J} of

$$\text{conv} \left(\left\{ \lim_{n \rightarrow +\infty} \partial f(\mathbf{v}_n) : (\mathbf{v}_n)_{n=1}^{+\infty} \text{ s.t. } \mathbf{v}_n \in \mathcal{E} \setminus \Omega, \mathbf{v}_n \xrightarrow{n \rightarrow +\infty} \mathbf{w} \right\} \right).$$

For a continuously differentiable function $f : \mathcal{E} \rightarrow \mathcal{F}$ or $f : \mathcal{E} \rightarrow \mathbb{R}$, there is a unique generalized gradient, recovering the usual gradient (Clarke *et al.*, 2008, Proposition 3.1, page 78). The chain rule can be generalized to these objects (Clarke *et al.*, 2008). Recently, Bolte and Pauwels (2020) and Bolte *et al.* (2022) further generalized Clarke gradients through the definition of conservative gradients to define automatic differentiation schemes for nonsmooth functions.

2.8 Summary

- The usual definition of **derivatives** of real-valued univariate functions extends to multivariate functions $f : \mathbb{R}^P \rightarrow \mathbb{R}$ through the notion of **directional derivative** $\partial f(\mathbf{w})[\mathbf{v}]$ at $\mathbf{w} \in \mathbb{R}^P$ in the direction $\mathbf{v} \in \mathbb{R}^P$.

- To take advantage of the representation of $\mathbf{w} = \sum_{j=1}^P w_j \mathbf{e}_j$ using the canonical bases $\{\mathbf{e}_1, \dots, \mathbf{e}_P\}$, the definition of **differentiable** functions requires the **linearity** of the directional derivative w.r.t. the direction \mathbf{v} .
- This requirement gives rise to the notion of **gradient** $\nabla f(\mathbf{w}) \in \mathbb{R}^P$, the vector that gathers the partial derivatives and further defines the **steepest ascent direction** at \mathbf{w} .
- For vector-input vector-output functions $f: \mathbb{R}^P \rightarrow \mathbb{R}^M$, the directional derivative leads to the definition of **Jacobian matrix** $\partial f(\mathbf{w}) \in \mathbb{R}^{M \times P}$, the matrix which gathers all partial derivatives (notice that we use bold ∂). The **chain rule** is then the **product** of Jacobian matrices.
- These notions can be extended to general Euclidean spaces, such as the spaces of matrices or tensors. For functions of the form $f: \mathcal{E} \rightarrow \mathbb{R}$, the gradient is $\nabla f(\mathbf{w}) \in \mathcal{E}$. More generally, for functions of the form $f: \mathcal{E} \rightarrow \mathcal{F}$, the Jacobian $\partial f(\mathbf{w})$ can be seen as a **linear map** (notice the non-bold ∂). The directional derivative at $\mathbf{w} \in \mathcal{E}$ naturally defines a linear map $l[\mathbf{v}] = \partial f(\mathbf{w})[\mathbf{v}]$, where $\partial f(\mathbf{w}): \mathcal{E} \rightarrow \mathcal{F}$ is called the **Jacobian vector product** (JVP) and captures the infinitesimal variation at $\mathbf{w} \in \mathcal{E}$ along the **input** direction $\mathbf{v} \in \mathcal{E}$.
- Its **adjoint** $\partial f(\mathbf{w})^*: \mathcal{F} \rightarrow \mathcal{E}$ defines another linear map $l[\mathbf{u}] = \partial f(\mathbf{w})^*[\mathbf{u}]$ called the **vector Jacobian product** (VJP) and captures the infinitesimal variation at $\mathbf{w} \in \mathcal{E}$ along the **output** direction $\mathbf{u} \in \mathcal{F}$. The **chain rule** is then the **composition** of these linear maps.
- For the particular case when we compose a scalar-valued function ℓ (such as a loss function) with a vector-valued function f (such as a network function), the gradient is given by $\nabla(\ell \circ f)(\mathbf{w}) = \partial f(\mathbf{w})^* \nabla \ell(f(\mathbf{w}))$. This is why being able to apply the adjoint to a gradient, which as we shall see can be done with reverse-mode autodiff, is so pervasive in machine learning.

- The definitions of JVP and VJP operators can further be generalized in the context of differentiable geometry. In that framework, the JVP amounts to the **pushforward** operator that acts on **tangent vectors**. The VJP amounts to the **pullback** operator that acts on **cotangent vectors**.
- We also saw that the **Hessian matrix** of a function $f(\mathbf{w})$ from \mathbb{R}^P to \mathbb{R} is denoted $\nabla^2 f(\mathbf{w}) \in \mathbb{R}^{P \times P}$. It is symmetric if the second partial derivatives are continuous. Seen as linear map, the Hessian leads to the notion of **Hessian-vector product** (HVP), which we saw can be reduced to the JVP or the VJP of $\nabla f(\mathbf{w})$.
- The main take-away message of this chapter is that computing the directional derivative or the gradient of compositions of functions **does not** require computing intermediate Jacobians but only to evaluate linear maps (JVPs or VJPs) associated with these intermediate functions. The goal of automatic differentiation, presented in Chapter 8, is precisely to provide an efficient implementation of these maps for **computation chains** or more generally for **computation graphs**.

3

Probabilistic learning

In this chapter, we review how to perform probabilistic learning. We also introduce exponential family distributions, as they play a key role in this book.

3.1 Probability distributions

3.1.1 Discrete probability distributions

A discrete probability distribution over a set \mathcal{Y} is specified by its **probability mass function** (PMF) $p: \mathcal{Y} \rightarrow [0, 1]$. The probability of $\mathbf{y} \in \mathcal{Y}$ is then defined by

$$\mathbb{P}(Y = \mathbf{y}) := p(\mathbf{y}),$$

where Y denotes a random variable. When Y follows a distribution p , we write $Y \sim p$ (with some abuse of notation, we use the same letter p to denote the distribution and the PMF). The **expectation** of $\phi(Y)$, where $Y \sim p$ and $\phi: \mathcal{Y} \rightarrow \mathbb{R}^M$, is then

$$\mathbb{E}[\phi(Y)] = \sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{y}) \phi(\mathbf{y}),$$

its **variance** (for one-dimensional variables) is

$$\mathbb{V}[\phi(Y)] = \mathbb{E}[(\phi(Y) - \mathbb{E}[\phi(Y)])^2] = \sum_{y \in \mathcal{Y}} p(y)(\phi(y) - \mathbb{E}[\phi(Y)])^2$$

and its **mode** is

$$\arg \max_{y \in \mathcal{Y}} p(y).$$

The **Kullback-Leibler** (KL) divergence (also known as relative entropy) between two discrete distributions over \mathcal{Y} , with associated PMFs p and q , is the statistical “distance” defined by

$$\text{KL}(p, q) := \sum_{y \in \mathcal{Y}} p(y) \log \left(\frac{p(y)}{q(y)} \right) = \mathbb{E}_{Y \sim p} \log \left(\frac{p(Y)}{q(Y)} \right).$$

3.1.2 Continuous probability distributions

A continuous probability distribution over \mathcal{Y} is specified by its **probability density function** (PDF) $p: \mathcal{Y} \rightarrow \mathbb{R}_+$. The probability of $\mathcal{A} \subseteq \mathcal{Y}$ is then

$$\mathbb{P}(Y \in \mathcal{A}) = \int_{\mathcal{A}} p(y) dy.$$

The definitions of expectation, variance and KL divergence are defined analogously to the discrete setting, simply replacing $\sum_{y \in \mathcal{Y}}$ with $\int_{\mathcal{Y}}$. Specifically, the expectation of $\phi(Y)$ is

$$\mathbb{E}[\phi(Y)] = \int_{\mathcal{Y}} p(y) \phi(y) dy,$$

the variance is

$$\mathbb{V}[\phi(Y)] = \mathbb{E}[(\phi(Y) - \mathbb{E}[\phi(Y)])^2] = \int_{\mathcal{Y}} p(y)(\phi(y) - \mathbb{E}[\phi(Y)])^2 dy$$

and the KL divergence is

$$\text{KL}(p, q) := \int_{\mathcal{Y}} p(y) \log \left(\frac{p(y)}{q(y)} \right) dy = \mathbb{E}_{Y \sim p} \log \left(\frac{p(Y)}{q(Y)} \right).$$

The mode is defined as the arg maximum of the PDF.

When $\mathcal{Y} = \mathbb{R}$, we can also define the **cumulative distribution function** (CDF)

$$\mathbb{P}(Y \leq b) = \int_{-\infty}^b p(y) dy.$$

The probability of Y lying in the semi-closed interval $(a, b]$ is then

$$\mathbb{P}(a < Y \leq b) = \mathbb{P}(Y \leq b) - \mathbb{P}(Y \leq a).$$

3.2 Maximum likelihood estimation

3.2.1 Negative log-likelihood

We saw that a probability distribution over \mathcal{Y} is specified by $p(\mathbf{y})$, which is called the probability mass function (PMF) for discrete variables or the probability density function (PDF) for continuous variables. In practice, the true distribution p generating the data is unknown and we wish to approximate it with a distribution p_{λ} , with parameters $\lambda \in \Lambda$. Given a finite set of i.i.d. observations $\mathbf{y}_1, \dots, \mathbf{y}_N$, how do we fit $\lambda \in \Lambda$ to the data? This can be done by maximizing the **likelihood** of the data, i.e., we seek to solve

$$\hat{\lambda}_N := \arg \max_{\lambda \in \Lambda} \prod_{i=1}^N p_{\lambda}(\mathbf{y}_i).$$

This is known as **maximum likelihood estimation** (MLE). Because the log function is monotonically increasing, this is equivalent to minimizing the **negative log-likelihood**, i.e., we have

$$\hat{\lambda}_N = \arg \min_{\lambda \in \Lambda} - \sum_{i=1}^N \log p_{\lambda}(\mathbf{y}_i).$$

Example 3.1 (MLE for the normal distribution). Suppose we set p_{λ} to the normal distribution with parameters $\lambda = (\mu, \sigma)$, i.e.,

$$p_{\lambda}(y) := \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{y-\mu}{\sigma}\right)^2\right).$$

Then, given observations y_1, \dots, y_N , the MLE estimators for μ and σ^2 are the sample mean and the sample variance, respectively.

3.2.2 Consistency w.r.t. the Kullback-Leibler divergence

It is well-known that the MLE estimator is **consistent**, in the sense of the Kullback-Leibler divergence. That is, denoting the true distribution

p and

$$\boldsymbol{\lambda}_\infty := \arg \min_{\boldsymbol{\lambda} \in \Lambda} \text{KL}(p, p_{\boldsymbol{\lambda}}) = \mathbb{E}_{Y \sim p} \log \left(\frac{p(Y)}{p_{\boldsymbol{\lambda}}(Y)} \right),$$

then $\hat{\boldsymbol{\lambda}}_N \rightarrow \boldsymbol{\lambda}_\infty$ in expectation over the observations, as $N \rightarrow \infty$. This can be seen by using

$$\text{KL}(p, p_{\boldsymbol{\lambda}}) \approx \frac{1}{N} \sum_{i=1}^N \log \left(\frac{p(\mathbf{y}_i)}{p_{\boldsymbol{\lambda}}(\mathbf{y}_i)} \right) = \frac{1}{N} \sum_{i=1}^N \log p(\mathbf{y}_i) - \log p_{\boldsymbol{\lambda}}(\mathbf{y}_i)$$

and the law of large numbers.

3.3 Probabilistic supervised learning

3.3.1 Conditional probability distributions

Many times in machine learning, instead of a probability $\mathbb{P}(Y = \mathbf{y})$ for some $\mathbf{y} \in \mathcal{Y}$, we wish to define a conditional probability $\mathbb{P}(Y = \mathbf{y} | X = \mathbf{x})$, for some input $\mathbf{x} \in \mathcal{X}$. This can be achieved by reduction to an unconditional probability distribution,

$$\mathbb{P}(Y = \mathbf{y} | X = \mathbf{x}) := p_{\boldsymbol{\lambda}}(\mathbf{y})$$

where

$$\boldsymbol{\lambda} := f(\mathbf{x}, \mathbf{w})$$

and f is a **model function** with **model parameters** $\mathbf{w} \in \mathcal{W}$. That is, rather than being a deterministic function from \mathcal{X} to \mathcal{Y} , f is a function from \mathcal{X} to Λ , the set of permissible **distribution parameters** of the output distribution associated with the input.

We emphasize that $\boldsymbol{\lambda}$ could be a single parameter or a collection of parameters. For instance, in the Bernoulli distribution, $\lambda = \pi$, while in the univariate normal distribution, $\boldsymbol{\lambda} = (\mu, \sigma)$.

In Section 3.4 and throughout this book, we will also use the notation $p_{\boldsymbol{\theta}}$ instead of $p_{\boldsymbol{\lambda}}$ when $\boldsymbol{\theta}$ are the canonical parameters of an exponential family distribution.

3.3.2 Inference

The main advantage of this probabilistic approach is that our prediction model is much richer than if we just learned a function from \mathcal{X} to \mathcal{Y} .

We now have access to the whole distribution over possible outcomes in \mathcal{Y} and can compute various statistics:

- Probability: $\mathbb{P}(Y = \mathbf{y} | X = \mathbf{x})$ or $\mathbb{P}(Y \in \mathcal{A} | X = \mathbf{x})$,
- Expectation: $\mathbb{E}[\phi(Y) | X = \mathbf{x}]$ for some function ϕ ,
- Variance: $\mathbb{V}[\phi(Y) | X = \mathbf{x}]$,
- Mode: $\arg \max_{\mathbf{y} \in \mathcal{Y}} p_{\lambda}(\mathbf{y})$.

We now review probability distributions useful for binary classification, multiclass classification, regression, multivariate regression, and integer regression. In the following, to make the notation more lightweight, we omit the dependence on \mathbf{x} .

3.3.3 Binary classification

For **binary outcomes**, where $\mathcal{Y} = \{0, 1\}$, we can use a **Bernoulli distribution** with parameter

$$\lambda := \pi \in [0, 1].$$

When a random variable Y is distributed according to a Bernoulli distribution with parameter π , we write

$$Y \sim \text{Bernoulli}(\pi).$$

The PMF of this distribution is

$$p_{\pi}(y) := \begin{cases} \pi & \text{if } y = 1 \\ 1 - \pi & \text{if } y = 0 \end{cases}.$$

The Bernoulli distribution is a **binomial distribution** with a single trial. Since $y \in \{0, 1\}$, the PMF can be rewritten as

$$p_{\pi}(y) = \pi^y (1 - \pi)^{1-y}.$$

The mean is

$$\mathbb{E}[Y] = \pi = \mathbb{P}(Y = 1)$$

and the variance is

$$\mathbb{V}[Y] = \pi(1 - \pi) = \mathbb{P}(Y = 1)\mathbb{P}(Y = 0).$$

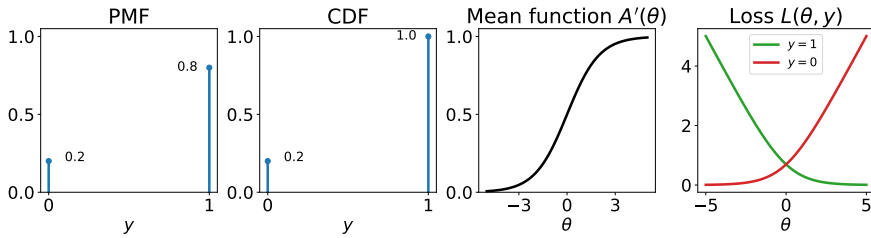


Figure 3.1: The **Bernoulli distribution**, whose PMF and CDF are here illustrated with parameter $\pi = 0.8$. Its mean function is $\pi = A'(\theta) = \text{logistic}(\theta) = \frac{1}{1 + \exp(-\theta)}$, where θ is for instance the output of a neural network. The negative log-likelihood leads to the **logistic loss**, $L(\theta, y) = \text{softplus}(\theta) - \theta y = \log(1 + \exp(\theta)) - \theta y$. The loss curve is shown for $y \in \{0, 1\}$.

Parameterization using a sigmoid

Since the parameter π of a Bernoulli distribution needs to belong to $[0, 1]$, we typically use a **sigmoid function** (Section 4.4.3), such as a **logistic function** as the output layer:

$$\pi := f(\mathbf{x}, \mathbf{w}) := \text{logistic}(g(\mathbf{x}, \mathbf{w})),$$

where $g: \mathcal{X} \times \mathcal{W} \rightarrow \mathbb{R}$ is for example a neural network and

$$\text{logistic}(a) := \frac{1}{1 + \exp(-a)} \in (0, 1).$$

When g is linear in \mathbf{w} , this is known as **binary logistic regression**.

Remark 3.1 (Link with the logistic distribution). The logistic distribution with mean and scale parameters μ and σ is a **continuous** probability distribution with PDF

$$p_{\mu, \sigma}(u) := p_{0,1}\left(\frac{u - \mu}{\sigma}\right)$$

where

$$p_{0,1}(z) := \frac{\exp(-z)}{(1 + \exp(-z))^2}.$$

If a random variable U follows a logistic distribution with parameters μ and σ , we write $U \sim \text{Logistic}(\mu, \sigma)$. The CDF of

$U \sim \text{Logistic}(\mu, \sigma)$ is

$$\mathbb{P}(U \leq u) = \int_{-\infty}^u p_{\mu, \sigma}(u) du = \text{logistic} \left(\frac{u - \mu}{\sigma} \right).$$

Therefore, if

$$U \sim \text{Logistic}(\mu, \sigma)$$

and

$$Y \sim \text{Bernoulli} \left(\text{logistic} \left(\frac{u - \mu}{\sigma} \right) \right),$$

then

$$\mathbb{P}(Y = 1) = \mathbb{P}(U \leq u).$$

Here, U can be interpreted as a latent continuous variable and u as a threshold.

3.3.4 Multiclass classification

For **categorical outcomes** with M possible choices, where $\mathcal{Y} = [M]$, we can use a **categorical distribution** with parameters

$$\boldsymbol{\lambda} := \boldsymbol{\pi} \in \Delta^M,$$

where we define the probability simplex

$$\Delta^M := \{\boldsymbol{\pi} \in \mathbb{R}_+^M : \langle \boldsymbol{\pi}, \mathbf{1} \rangle = 1\},$$

i.e., the set of valid discrete probability distributions. When Y follows a categorical distribution with parameter $\boldsymbol{\pi}$, we write

$$Y \sim \text{Categorical}(\boldsymbol{\pi}).$$

The PMF of the categorical distribution is

$$p_{\boldsymbol{\pi}}(y) := \langle \boldsymbol{\pi}, \phi(y) \rangle = \pi_y,$$

where

$$\phi(y) := \mathbf{e}_y$$

is the standard basis vector for the coordinate $y \in [M]$.

Since Y is a categorical variable, it does not make sense to compute the expectation of Y but we can compute that of $\phi(Y) = \mathbf{e}_Y$,

$$\mathbb{E}_{Y \sim p_{\boldsymbol{\pi}}}[\phi(Y)] = \boldsymbol{\pi}.$$

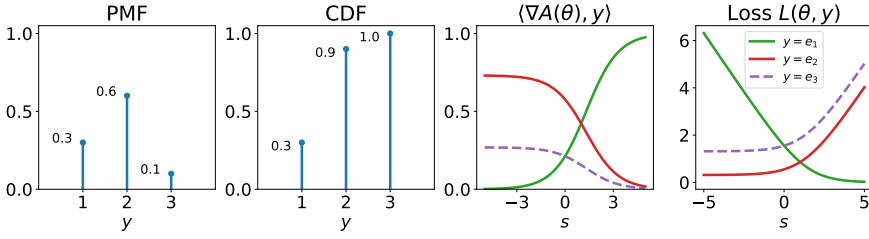


Figure 3.2: The **categorical distribution**, whose PMF and CDF are here illustrated with parameter $\pi = (0.3, 0.6, 0.1)$. Its mean function is $\pi = \nabla A(\theta) = \text{softargmax}(\theta)$, where $\theta \in \mathbb{R}^M$ is for instance the output of a neural network. Here, for illustration purpose, we choose to set $\theta = (s, 1, 0)$ and vary only s . Since the mean function $\nabla A(\theta)$ belongs to \mathbb{R}^3 , we choose to display $\langle \nabla A(\theta), e_i \rangle = \nabla A(\theta)_i$, for $i \in \{1, 2, 3\}$. The negative log-likelihood leads to the **logistic loss**, $L(\theta, y) = \text{logsumexp}(\theta) - \langle \theta, y \rangle$. The loss curve is shown for $y \in \{e_1, e_2, e_3\}$, again with $\theta = (s, 1, 0)$ and varying s .

Therefore, as was also the case for the Bernoulli distribution, the mean and the probability distribution (represented by the vector π) are the same in this case.

Parameterization using a softargmax

Since the parameter vector π of a categorical distribution needs to belong to Δ^M , we typically use a softargmax as the output layer:

$$\pi := f(x, w) := \text{softargmax}(g(x, w)),$$

where $g: \mathcal{X} \times \mathcal{W} \rightarrow \mathbb{R}^M$ is for example a neural network and

$$\text{softargmax}(u) := \frac{\exp(u)}{\sum_j \exp(u_j)} \in \text{relint}(\Delta^M).$$

The output of the softargmax is in the relative interior of Δ^M , $\text{relint}(\Delta^M) = \Delta^M \cap \mathbb{R}_{>0}^M$. That is, the produced probabilities are always strictly positive. The categorical distribution is a **multinomial distribution** with a single trial. When g is linear in w , this is therefore known as **multi-class** or **multinomial logistic regression**, though strictly speaking a multinomial distribution could use more than one trial.

3.3.5 Regression

For **real outcomes**, where $\mathcal{Y} = \mathbb{R}$, we can use, among other choices, a **normal distribution** with parameters

$$\boldsymbol{\lambda} := (\mu, \sigma),$$

where $\mu \in \mathbb{R}$ is the mean parameter and $\sigma \in \mathbb{R}_+$ is the standard deviation parameter. When Y follows a normal distribution with parameters (μ, σ) , we write

$$Y \sim \text{Normal}(\mu, \sigma).$$

The PDF is

$$p_{\mu, \sigma}(y) := \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2} \frac{(y - \mu)^2}{\sigma^2}\right).$$

The expectation is

$$\mathbb{E}_{Y \sim p_{\mu, \sigma}}[Y] = \mu.$$

One advantage of the probabilistic perspective is that we are not limited to predicting the mean. We can also compute the CDF

$$\mathbb{P}(Y \leq y) = \frac{1}{2} \left[1 + \operatorname{erf}\left(\frac{y - \mu}{\sigma\sqrt{2}}\right) \right],$$

where we used the **error function**

$$\operatorname{erf}(z) := \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt.$$

This function is available in most scientific computing libraries, such as SciPy (Virtanen *et al.*, 2020). We can also write

$$\mathbb{P}(Y \leq y) = \Phi\left(\frac{y - \mu}{\sigma}\right)$$

where

$$\Phi(z) := \frac{1}{2} \left[1 + \operatorname{erf}\left(\frac{z}{\sqrt{2}}\right) \right] \quad (3.1)$$

is the CDF of the standard Gaussian distribution (with zero mean and unit variance). From the CDF, we also easily obtain

$$\mathbb{P}(a < Y \leq b) = \frac{1}{2} \left[\operatorname{erf}\left(\frac{b - \mu}{\sigma\sqrt{2}}\right) - \operatorname{erf}\left(\frac{a - \mu}{\sigma\sqrt{2}}\right) \right].$$

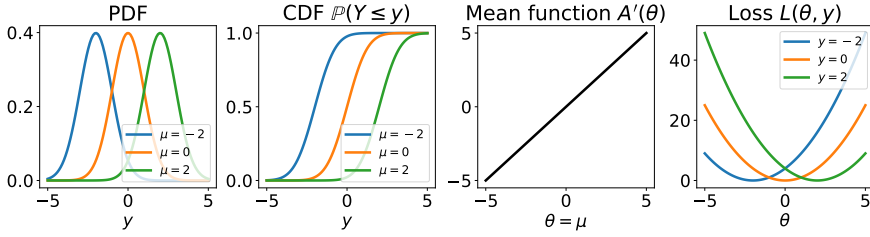


Figure 3.3: The **Gaussian distribution**, with mean parameter μ and variance $\sigma^2 = 1$. Its mean function is $\mu = A'(\theta) = \theta$, where θ is for instance the output of a neural network. The negative log-likelihood leads to the **squared loss**, $L(\theta, y) = (y - \theta)^2$. The loss curve is shown for $y \in \{-2, 0, 2\}$.

Parameterization

Typically, in regression, the mean is output by a model, while the standard deviation σ is kept fixed (typically set to 1). Since μ is unconstrained, we can simply set

$$\mu := f(\mathbf{x}, \mathbf{w}) \in \mathbb{R},$$

where $f: \mathcal{X} \times \mathcal{W} \rightarrow \mathbb{R}$ is for example a neural network. That is, the output of f is the mean of the distribution,

$$\mathbb{E}_{Y \sim p_{\mu,1}}[Y] = \mu = f(\mathbf{x}, \mathbf{w}).$$

We can also use μ to predict $\mathbb{P}(Y \leq y)$ or $\mathbb{P}(a < Y \leq b)$, as shown above.

3.3.6 Multivariate regression

More generally, for **multivariate outcomes**, where $\mathcal{Y} = \mathbb{R}^M$, we can use a **multivariate normal distribution** with parameters

$$\lambda := (\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

where $\boldsymbol{\mu} \in \mathbb{R}^M$ is the mean and $\boldsymbol{\Sigma} \in \mathbb{R}^{M \times M}$ is the covariance matrix. When Y follows a multivariate normal distribution with parameters $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, we write

$$Y \sim \text{Normal}(\boldsymbol{\mu}, \boldsymbol{\sigma}).$$

The PDF is

$$p_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}(\mathbf{y}) := \frac{1}{\sqrt{2\pi^M |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2} \langle \mathbf{y} - \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \boldsymbol{\mu}) \rangle\right).$$

Using a diagonal covariance matrix is equivalent to using M independent normal distributions for each Y_j , for $j \in [M]$. The expectation is

$$\mathbb{E}_{Y \sim p_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}}[Y] = \boldsymbol{\mu}.$$

Parameterization

Typically, in multivariate regression, the mean is output by a model, while the covariance matrix is kept fixed (typically set to the identity matrix). Since $\boldsymbol{\mu}$ is again unconstrained, we can simply set

$$\boldsymbol{\mu} := f(\mathbf{x}, \mathbf{w}) \in \mathbb{R}^M.$$

More generally, we can parametrize the function f so as to output both the mean $\boldsymbol{\mu}$ and the covariance matrix $\boldsymbol{\Sigma}$, i.e.,

$$(\boldsymbol{\mu}, \boldsymbol{\Sigma}) := f(\mathbf{x}, \mathbf{w}) \in \mathbb{R}^M \times \mathbb{R}^{M \times M}.$$

The function f must be designed such that $\boldsymbol{\Sigma}$ is symmetric and positive semi-definite. This is easy to achieve for instance by parametrizing $\boldsymbol{\Sigma} = \mathbf{S}\mathbf{S}^\top$ for some matrix \mathbf{S} .

3.3.7 Integer regression

For **integer outcomes**, where $\mathcal{Y} = \mathbb{N}$, we can use, among other choices, a **Poisson distribution** with mean parameter $\lambda > 0$. When Y follows a Poisson distribution with parameter λ , we write

$$Y \sim \text{Poisson}(\lambda).$$

The PMF is

$$\mathbb{P}(Y = y) = p_\lambda(y) := \frac{\lambda^y \exp(-\lambda)}{y!}.$$

It is the probability of y events occurring in an interval of time. The Poisson distribution is frequently used when there is a large number of possible events, each of which is rare.

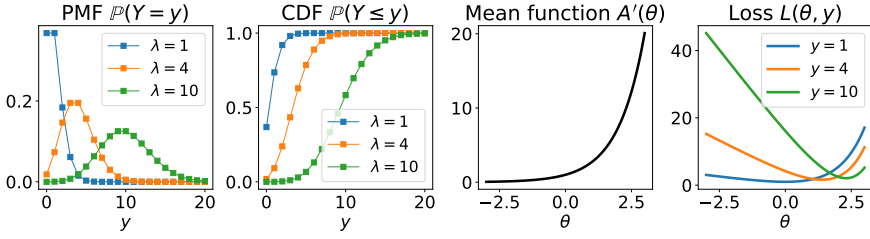


Figure 3.4: The **Poisson distribution**, with mean parameter λ . For the PMF and the CDF, the lines between markers are shown for visual aid: the Poisson distribution does not assign probability mass to non-integer values. Its mean function is $\lambda = A'(\theta) = \exp(\theta)$, where θ is for instance the output of a neural network. The negative log-likelihood leads to the **Poisson loss**, $L(\theta, y) = -\log p_\lambda(y) = -y\theta + \exp(\theta) + \log(y!)$, which is a convex function of θ . The loss curve is shown for $y \in \{1, 4, 10\}$.

The CDF is

$$\mathbb{P}(Y \leq y) = \sum_{i=0}^y \mathbb{P}(Y = i).$$

The Poisson distribution implies that the **index of dispersion** (the ratio between variance and mean) is 1, since

$$\mathbb{E}[Y] = \mathbb{V}[Y] = \lambda.$$

When this assumption is inappropriate, one can use generalized Poisson distributions (Satterthwaite, 1942).

Parameterization using an exponential

Since the parameter λ of a Poisson distribution needs to be strictly positive, we typically use an exponential function as output layer:

$$\lambda := f(\mathbf{x}, \mathbf{w}) := \exp(g(\mathbf{x}, \mathbf{w})) > 0,$$

where $g: \mathcal{X} \times \mathcal{W} \rightarrow \mathbb{R}$.

3.3.8 Loss functions

We now discuss how to learn the model parameters $\mathbf{w} \in \mathcal{W}$ from input-output pairs $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)$.

Deterministic vs. probabilistic approaches

In a deterministic approach, if we used a mapping $f: \mathcal{X} \times \mathcal{W} \rightarrow \mathcal{Y}$, we could formulate an objective function of the form

$$L(\mathbf{w}) := \frac{1}{N} \sum_{i=1}^N \ell(f(\mathbf{x}_i, \mathbf{w}), \mathbf{y}_i),$$

where $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a loss function. Unfortunately, $L(\mathbf{w})$ would be typically **discontinuous** if \mathcal{Y} is a discrete output space (as is the case in classification), making optimization difficult.

In contrast, in the probabilistic approach, we use a mapping $f: \mathcal{X} \times \mathcal{W} \rightarrow \Lambda$ to distribution parameters, and we can formulate an objective of the form

$$L(\mathbf{w}) := \frac{1}{N} \sum_{i=1}^N \ell(f(\mathbf{x}_i, \mathbf{w}), \mathbf{y}_i),$$

where $\ell: \Lambda \times \mathcal{Y} \rightarrow \mathbb{R}$. This is typically a **continuous** objective, since Λ is typically a continuous set and p_{λ} varies continuously w.r.t. λ even if p_{λ} is a distribution over a discrete set \mathcal{Y} . In other words, the probabilistic approach is not only powerful for the inference it allows us to do (probability, expectation, variance, mode), but also because it allows to formulate a continuous and typically differentiable objective function!

Negative log-likelihood

In the conditional setting briefly reviewed in Section 3.3.1, we can use maximum likelihood estimation (MLE) to estimate the model parameters $\mathbf{w} \in \mathcal{W}$ of f . Given a set of input-output pairs $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)$, we choose the model parameters that maximize the **likelihood** of the data,

$$\hat{\mathbf{w}}_N := \arg \max_{\mathbf{w} \in \mathcal{W}} \prod_{i=1}^N p_{\lambda_i}(\mathbf{y}_i),$$

where

$$\lambda_i := f(\mathbf{x}_i, \mathbf{w}).$$

Again, this is equivalent to minimizing the **negative log-likelihood**,

$$\hat{\mathbf{w}}_N = \arg \min_{\mathbf{w} \in \mathcal{W}} - \sum_{i=1}^N \log p_{\lambda_i}(\mathbf{y}_i).$$

In the notation above, this corresponds to defining the loss function

$$\ell(\lambda_i, \mathbf{y}_i) := -\log p_{\lambda_i}(\mathbf{y}_i).$$

Recovering well-known loss functions

Interestingly, MLE allows us to recover several popular loss functions.

- For the Bernoulli distribution with parameter $\lambda_i = \pi_i = \text{logistic}(g(\mathbf{x}_i, \mathbf{w}))$, we have

$$-\log p_{\lambda_i}(y_i) = -[y_i \log \pi_i + (1 - y_i) \log(1 - \pi_i)],$$

which is the **binary logistic loss** function.

- For the categorical distribution with parameters $\lambda_i = \pi_i = \text{softargmax}(g(\mathbf{x}_i, \mathbf{w}))$, we have

$$\begin{aligned} -\log p_{\lambda_i}(y_i) &= \log \sum_{j=1}^M \exp(\pi_{i,j} - \pi_{i,y_i}) \\ &= \text{logsumexp}(\boldsymbol{\pi}_i) - \langle \boldsymbol{\pi}_i, \mathbf{e}_{y_i} \rangle, \end{aligned}$$

which is the **multiclass logistic loss** function, also known as **cross-entropy loss**.

- For the normal distribution with mean $\lambda_i = \mu_i = f(\mathbf{x}_i, \mathbf{w})$ and fixed variance σ_i^2 , we have

$$-\log p_{\lambda_i}(y_i) = \frac{1}{\sigma_i^2} (y_i - \mu_i)^2 + \frac{1}{2} \log \sigma_i^2 + \frac{1}{2} \log(2\pi),$$

which is, up to constant and with unit variance, the **squared loss** function.

- For the Poisson distribution with mean $\lambda_i = \exp(\theta_i)$, where $\theta_i := g(\mathbf{x}_i, \mathbf{w})$, we have

$$\begin{aligned} -\log p_{\lambda_i}(y_i) &= -y_i \log(\lambda_i) + \lambda_i + \log(y_i!) \\ &= -y_i \theta_i + \exp(\theta_i) + \log(y_i!) \end{aligned}$$

which is the **Poisson loss** function. The loss function is convex w.r.t. λ_i and θ_i for $y_i \geq 0$.

3.4 Exponential family distributions

3.4.1 Definition

The exponential family is a class of probability distributions, whose PMF or PDF can be written in the form

$$\begin{aligned} p_{\theta}(\mathbf{y}) &= \frac{h(\mathbf{y}) \exp[\langle \theta, \phi(\mathbf{y}) \rangle]}{\exp(A(\theta))} \\ &= h(\mathbf{y}) \exp[\langle \theta, \phi(\mathbf{y}) \rangle - A(\theta)], \end{aligned}$$

where θ are the **natural** or **canonical parameters** of the distribution. The function h is known as the base measure. The function ϕ is the **sufficient statistic**: it holds all the information about \mathbf{y} and is used to embed \mathbf{y} in a vector space. The function A is the **log-partition** or **log-normalizer** (see below for a details). All the distributions we reviewed in Section 3.3 belong to the exponential family. With some abuse of notation, we use p_{λ} for the distribution in **original form** and p_{θ} for the distribution in **exponential family form**. As we will see, we can go from θ to λ and vice-versa. We illustrate how to rewrite a distribution in exponential family form below.

Example 3.2 (Bernoulli distribution). The PMF of the Bernoulli distribution with parameter $\lambda = \pi$ equals

$$\begin{aligned} p_{\lambda}(y) &:= \pi^y (1 - \pi)^{1-y} \\ &= \exp(\log(\pi^y (1 - \pi)^{1-y})) \\ &= \exp(y \log(\pi) + (1 - y) \log(1 - \pi)) \\ &= \exp(\log(\pi / (1 - \pi))y + \log(1 - \pi)) \\ &= \exp(\theta y - \log(1 + \exp(\theta))) \\ &= \exp(\theta y - \text{softplus}(\theta)) \\ &=: p_{\theta}(y). \end{aligned}$$

Therefore, Bernoulli distributions belong to the exponential family,

with natural parameter $\theta = \text{logit}(\pi) := \log(\pi/(1 + \pi))$. Conversely, we have $\pi = \text{logistic}(\theta) = \frac{1}{1 + \exp(-\theta)}$.

We rewrite the previously-described distributions in exponential family form in Table 3.1. This list is non-exhaustive: there are many more distributions in the exponential family! (Barndorff-Nielsen, 2014)

3.4.2 The log-partition function

The log-partition function A is the logarithm of the distribution's normalization factor. That is,

$$A(\boldsymbol{\theta}) := \log \sum_{\mathbf{y} \in \mathcal{Y}} h(\mathbf{y}) \exp [\langle \boldsymbol{\theta}, \phi(\mathbf{y}) \rangle]$$

for discrete random variables and

$$A(\boldsymbol{\theta}) := \log \int_{\mathcal{Y}} h(y) \exp [\langle \boldsymbol{\theta}, \phi(y) \rangle] dy$$

for continuous random variables. We denote the set of valid parameters

$$\Theta := \{\boldsymbol{\theta} \in \mathbb{R}^M : A(\boldsymbol{\theta}) < +\infty\} \subseteq \mathbb{R}^M.$$

We can conveniently rewrite $A(\boldsymbol{\theta})$ for discrete random variables as

$$A(\boldsymbol{\theta}) = \text{logsumexp}(B(\boldsymbol{\theta})) := \log \sum_{\mathbf{y} \in \mathcal{Y}} [B(\boldsymbol{\theta})]_{\mathbf{y}},$$

and similarly for continuous variables. Here, we defined the **affine map**

$$B(\boldsymbol{\theta}) := (\langle \boldsymbol{\theta}, \phi(\mathbf{y}) \rangle + \log h(\mathbf{y}))_{\mathbf{y} \in \mathcal{Y}}.$$

Since $A(\boldsymbol{\theta})$ is the composition of logsumexp , a convex function, and of B , an affine map, we immediately obtain the following proposition.

Proposition 3.1 (Convexity of the log-partition). $A(\boldsymbol{\theta})$ is a convex function.

A major property of the log-partition function is that its gradient coincides with the expectation of $\phi(Y)$ according to $p_{\boldsymbol{\theta}}$.

Table 3.1: Examples of distributions in the exponential family.

	Bernoulli	Categorical
\mathcal{Y}	$\{0, 1\}$	$[M]$
λ	$\pi = \text{logistic}(\theta)$	$\boldsymbol{\pi} = \text{softargmax}(\boldsymbol{\theta})$
$\boldsymbol{\theta}$	$\text{logit}(\pi)$	$\log \boldsymbol{\pi} + \exp(A(\boldsymbol{\theta}))$
$\phi(y)$	y	\mathbf{e}_y
$A(\boldsymbol{\theta})$	$\text{softplus}(\theta)$	$\text{logsumexp}(\boldsymbol{\theta})$
$h(y)$	1	1
	Normal (location only)	Normal (location-scale)
\mathcal{Y}	\mathbb{R}	\mathbb{R}
λ	$\mu = \theta\sigma$	$(\mu, \sigma^2) = (\frac{-\theta_1}{2\theta_2}, \frac{-1}{2\theta_2})$
$\boldsymbol{\theta}$	$\frac{\mu}{\sigma}$	$(\frac{\mu}{\sigma^2}, \frac{-1}{2\sigma^2})$
$\phi(y)$	$\frac{y}{\sigma}$	(y, y^2)
$A(\boldsymbol{\theta})$	$\frac{\theta^2}{2} = \frac{\mu^2}{2\sigma^2}$	$\frac{-\theta_1^2}{4\theta_2} - \frac{1}{2} \log(-2\theta_2) = \frac{\mu^2}{\sigma^2} + \log \sigma$
$h(y)$	$\frac{\exp(\frac{-y^2}{2\sigma^2})}{\sqrt{2\pi}\sigma}$	$\frac{1}{\sqrt{2\pi}}$
	Multivariate normal	Poisson
\mathcal{Y}	\mathbb{R}^M	\mathbb{N}
λ	$(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (-\frac{1}{2}\boldsymbol{\theta}_2^{-1}\boldsymbol{\theta}_1, -\frac{1}{2}\boldsymbol{\theta}_2^{-1})$	$\lambda = \exp(\theta)$
$\boldsymbol{\theta}$	$(\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}, -\frac{1}{2}\boldsymbol{\Sigma}^{-1})$	$\log \lambda$
$\phi(\mathbf{y})$	$(\mathbf{y}, \mathbf{y}\mathbf{y}^\top)$	y
$A(\boldsymbol{\theta})$	$-\frac{1}{4}\boldsymbol{\theta}_1^\top \boldsymbol{\theta}_2^{-1} - \frac{1}{2} \log -2\boldsymbol{\theta}_2 $ $= \frac{1}{2}\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \frac{1}{2} \log \boldsymbol{\Sigma} $	$\exp(\theta)$
$h(y)$	$(2\pi)^{-M/2}$	$1/y!$

Proposition 3.2 (Gradient of the log-partition).

$$\boldsymbol{\mu}(\boldsymbol{\theta}) := \nabla A(\boldsymbol{\theta}) = \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}}[\phi(Y)] \in \mathcal{M}.$$

Proof. The result follows directly from

$$\nabla A(\boldsymbol{\theta}) = \partial B(\boldsymbol{\theta})^* \nabla \log \text{sumexp}(B(\boldsymbol{\theta})) = (\phi(\mathbf{y}))_{\mathbf{y} \in \mathcal{Y}} \text{softmax}(B(\boldsymbol{\theta})).$$

□

The gradient $\nabla A(\boldsymbol{\theta})$ is therefore often called the **mean function**. The set of achievable means $\boldsymbol{\mu}(\boldsymbol{\theta})$ is defined by

$$\mathcal{M} := \text{conv}(\phi(\mathcal{Y})) := \{\mathbb{E}_p[\phi(Y)] : p \in \mathcal{P}(\mathcal{Y})\},$$

where $\text{conv}(\mathcal{S})$ is the convex hull of \mathcal{S} and $\mathcal{P}(\mathcal{Y})$ is the set of valid probability distributions over \mathcal{Y} .

Similarly, the Hessian $\nabla^2 A(\boldsymbol{\theta})$ coincides with the covariance matrix of $\phi(Y)$ according to $p_{\boldsymbol{\theta}}$ (Wainwright and Jordan, 2008, Chapter 3).

When the exponential family is **minimal**, which means that the parameters $\boldsymbol{\theta}$ uniquely identify the distribution, it is known that ∇A is a one-to-one mapping from Θ to \mathcal{M} . That is, $\boldsymbol{\mu}(\boldsymbol{\theta}) = \nabla A(\boldsymbol{\theta})$ and $\boldsymbol{\theta} = (\nabla A)^{-1}(\boldsymbol{\mu}(\boldsymbol{\theta}))$.

3.4.3 Maximum entropy principle

Suppose we observe the empirical mean $\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{i=1}^n \phi(\mathbf{y}_i) \in \mathcal{M}$ of some observations $\mathbf{y}_1, \dots, \mathbf{y}_N$. How do we find a probability distribution achieving this mean? Clearly, such a distribution may not be unique. One way to choose among all possible distributions is by using the **maximum entropy principle**. Let us define the Shannon entropy by

$$H(p) := - \sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{y}) \log p(\mathbf{y})$$

for discrete variables and by

$$H(p) := - \int_{\mathcal{Y}} p(y) \log p(y) dy$$

for continuous variables. This captures the level of “uncertainty” in p , i.e., it is maximized when the distribution is uniform. Then, the **maximum entropy distribution** satisfying the first-order moment condition (i.e., whose expectation matches the empirical mean) is

$$p^* := \arg \max_{p \in \mathcal{P}(\mathcal{Y})} H(p) \quad \text{s.t.} \quad \mathbb{E}_{Y \sim p}[\phi(Y)] = \hat{\boldsymbol{\mu}}.$$

It can be shown that the maximum entropy distribution is necessarily in the exponential family with sufficient statistics defined by ϕ and its canonical parameters $\boldsymbol{\theta}$ coincide with the **Lagrange multipliers** of the above constraint (Wainwright and Jordan, 2008, Section 3.1).

3.4.4 Maximum likelihood estimation

Similarly as in Section 3.2, to fit the parameters $\boldsymbol{\theta} \in \Theta$ of an exponential family distribution to some i.i.d. observations $\mathbf{y}_1, \dots, \mathbf{y}_N$, we can use the MLE principle, i.e.,

$$\hat{\boldsymbol{\theta}}_N = \arg \max_{\boldsymbol{\theta} \in \Theta} \prod_{i=1}^N p_{\boldsymbol{\theta}}(\mathbf{y}_i) = \arg \min_{\boldsymbol{\theta} \in \Theta} - \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{y}_i).$$

Fortunately, for exponential family distributions, the log probability/density enjoys a particularly simple form.

Proposition 3.3 (Negative log-likelihood). The negative log-likelihood of an exponential family distribution is

$$-\log p_{\boldsymbol{\theta}}(\mathbf{y}) = A(\boldsymbol{\theta}) - \langle \boldsymbol{\theta}, \phi(\mathbf{y}) \rangle - \log h(\mathbf{y}).$$

Its gradient is

$$-\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{y}) = \nabla A(\boldsymbol{\theta}) - \phi(\mathbf{y}) = \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}}[\phi(Y)] - \phi(\mathbf{y})$$

and its Hessian is

$$-\nabla_{\boldsymbol{\theta}}^2 \log p_{\boldsymbol{\theta}}(\mathbf{y}) = \nabla^2 A(\boldsymbol{\theta}),$$

which is **independent** of \mathbf{y} .

It follows from Proposition 3.1 that $\boldsymbol{\theta} \mapsto -\log p_{\boldsymbol{\theta}}(\mathbf{y})$ is **convex**. Interestingly, we see that the gradient is the **residual** between the

expectation of $\phi(Y)$ according to the model and the observed $\phi(\mathbf{y})$. Therefore, the negative log-likelihood of an exponential family distribution can be seen as performing first moment matching.

3.4.5 Probabilistic learning with exponential families

In the supervised probabilistic learning setting, we wish to estimate a conditional distribution of the form $p_{\mathbf{w}}(\mathbf{y} \mid \mathbf{x})$. Given a model function f , such as a neural network, a common approach for defining such a conditional distribution is by reduction to the unconditional setting,

$$p_{\mathbf{w}}(\mathbf{y} \mid \mathbf{x}) := p_{\boldsymbol{\theta}}(\mathbf{y}) \quad \text{where} \quad \boldsymbol{\theta} := f(\mathbf{x}, \mathbf{w}).$$

In other words, the role of f is to produce the parameters of $p_{\boldsymbol{\theta}}$ given some input \mathbf{x} . It is a function from $\mathcal{X} \times \mathcal{W}$ to Θ . Note that f must be designed such that it produces an output in

$$\Theta := \{\boldsymbol{\theta} \in \mathbb{R}^M : A(\boldsymbol{\theta}) < +\infty\}.$$

Many times, Θ will be the entire \mathbb{R}^M but this is not always the case. For instance, as we previously discussed, for a multivariate normal distribution, where $\boldsymbol{\theta} = (\boldsymbol{\mu}, \boldsymbol{\Sigma}) = f(\mathbf{x}, \mathbf{w})$, we need to ensure that $\boldsymbol{\Sigma}$ is a positive semidefinite matrix.

Training

Given input-output pairs $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$, we then seek to find the parameters \mathbf{w} of $f(\mathbf{x}, \mathbf{w})$ by minimizing the negative log-likelihood

$$\arg \min_{\mathbf{w} \in \mathcal{W}} - \sum_{i=1}^N \log p_{\boldsymbol{\theta}_i}(\mathbf{y}_i) = \arg \min_{\mathbf{w} \in \mathcal{W}} \sum_{i=1}^N A(\boldsymbol{\theta}_i) - \langle \boldsymbol{\theta}_i, \phi(\mathbf{y}_i) \rangle$$

where $\boldsymbol{\theta}_i := f(\mathbf{x}_i, \mathbf{w})$. While $-\log p_{\boldsymbol{\theta}}(\mathbf{y})$ is a convex function of $\boldsymbol{\theta}$ for exponential family distributions, we emphasize that $-\log p_{f(\mathbf{x}, \mathbf{w})}(\mathbf{y})$ is typically a nonconvex function of \mathbf{w} , when f is a **nonlinear** function, such as a neural network.

Inference

Once we found \mathbf{w} by minimizing the objective function above, there are several possible strategies to perform inference for a new input \mathbf{x} .

- **Expectation.** When the goal is to compute the expectation of $\phi(Y)$, we can use $\nabla A(f(\mathbf{x}, \mathbf{w}))$. That is, we compute the distribution parameters associated with \mathbf{x} by $\boldsymbol{\theta} = f(\mathbf{x}, \mathbf{w})$ and then we compute the mean by $\boldsymbol{\mu} = \nabla A(\boldsymbol{\theta})$. When f is linear in \mathbf{w} , the composition $\nabla A \circ f$ is called a **generalized linear model**.
- **Probability.** When the goal is to compute the probability of a certain \mathbf{y} , we can compute the distribution parameters associated with \mathbf{x} by $\boldsymbol{\theta} = f(\mathbf{x}, \mathbf{w})$ and then we can compute $\mathbb{P}(Y = \mathbf{y} | X = \mathbf{x}) = p_{\boldsymbol{\theta}}(\mathbf{y})$. In the particular case of the categorical distribution (of which the Bernoulli distribution is a special case), we point out again that the mean and the probability vector coincide:

$$\boldsymbol{\mu} = \mathbf{p} = \nabla A(\boldsymbol{\theta}) = \text{softargmax}(\boldsymbol{\theta}) \in \triangle^M.$$

- **Other statistics.** When the goal is to compute other quantities, such as the variance or the CDF, we can convert the natural parameters $\boldsymbol{\theta}$ to the original distribution parameters $\boldsymbol{\lambda}$ (see Table 3.1 for examples). Then, we can use established formulas for the distribution in original form, to compute the desired quantities.

3.5 Summary

- We reviewed **discrete** and **continuous** probability distributions.
- We saw how to fit distribution parameters to data using the **maximum likelihood estimation** (MLE) principle and saw its connection with the **Kullback-Leibler divergence**.
- Instead of designing a model function from the input space \mathcal{X} to the output space \mathcal{Y} , we saw that we can perform **probabilistic supervised learning** by designing a model function from \mathcal{X} to **distribution parameters** Λ .
- Leveraging the so-obtained parametric **conditional distribution** then allowed us to compute, not only output probabilities, but also various statistics such as the mean and the variance of the outputs.

- We reviewed the **exponential family**, a principled generalization of numerous distributions, which we saw is tightly connected with the **maximum entropy principle**.
- Importantly, the approaches described in this chapter produce perfectly valid **computation graphs**, meaning that we can combine them with neural networks and we can use automatic differentiation, to compute their derivatives.

Part II

Differentiable programs

4

Parameterized programs

Neural networks can be thought of as parameterized programs: programs with learnable parameters. In this chapter, we begin by reviewing how to represent programs mathematically. We then review several key neural network architectures and components.

4.1 Representing computer programs

4.1.1 Computation chains

To begin with, we consider simple programs that apply a **sequence** of functions f_1, \dots, f_K to an input $s_0 \in \mathcal{S}_0$. We call such programs **computation chains**. For example, an image may go through a sequence of transformations such as cropping, rotation, normalization, and so on. In neural networks, the transformations are typically parameterized, and the parameters are learned, leading to feedforward networks, presented in Section 4.2. Another example of sequence of functions is a for loop, presented in Section 5.8.

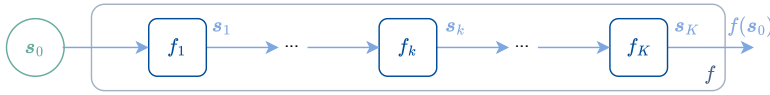


Figure 4.1: A computation chain is a sequence of function compositions. In the graph above, each intermediate node represents a single function. The first node represents the input, the last node the output. Edges represent the dependencies of the functions with respect to previous outputs or to the initial input.

Formally, a computation chain can be written as

$$\begin{aligned}
 s_0 &\in \mathcal{S}_0 \\
 s_1 &:= f_1(s_0) \in \mathcal{S}_1 \\
 &\vdots \\
 s_K &:= f_K(s_{K-1}) \in \mathcal{S}_K \\
 f(s_0) &:= s_K.
 \end{aligned} \tag{4.1}$$

Here, s_0 is the **input**, $s_k \in \mathcal{S}_k$ is an intermediate **state** of the program, and $s_K \in \mathcal{S}_K$ is the final **output**. Of course, the domain (input space) of f_k must be compatible with the image (output space) of f_{k-1} . That is, we should have $f_k: \mathcal{S}_{k-1} \rightarrow \mathcal{S}_k$. We can write a computation chain equivalently as

$$\begin{aligned}
 f(s_0) &= (f_K \circ \dots \circ f_2 \circ f_1)(s_0) \\
 &= f_K(\dots f_2(f_1(s_0))).
 \end{aligned}$$

A computation chain can be represented by a directed graph, shown in Fig. 4.1. The edges in the chain define a **total order**. The order is total, since two nodes are necessarily linked to each other by a path.

4.1.2 Directed acyclic graphs

In generic programs, intermediate functions may depend, not only on the previous function output, but on the outputs of several different functions. Such dependencies are best expressed using graphs.

A **directed graph** $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is defined by a set of **vertices** or **nodes** \mathcal{V} and a set of **edges** \mathcal{E} defining directed dependencies between

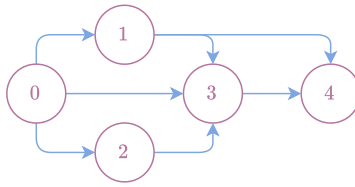


Figure 4.2: Example of a **directed acyclic graph**. Here the nodes are $\mathcal{V} = \{0, 1, 2, 3, 4\}$, the edges are $\mathcal{E} = \{(0, 1), (0, 2), (0, 3), (1, 3), (2, 3), (1, 4), (3, 4)\}$. Parents of the node 3 are $\text{pa}(3) = \{0, 1, 2\}$. Children of node 1 are $\text{ch}(1) = \{3, 4\}$. There is a unique root, 0, and a unique leaf, 4; $0 \rightarrow 3 \rightarrow 4$ is a path from 0 to 4. This is an acyclic graph since there is no cycle (i.e., a path from a node to itself). We can order nodes 0 and 3 as $0 \leq 3$ since there is no path from 3 to 0. Similarly, we can order 1 and 2 as $1 \leq 2$ since there is no path from 2 to 1. Two possible topological orders of the nodes are $(0, 1, 2, 3, 4)$ and $(0, 2, 1, 3, 4)$.

vertices. An edge $(i, j) \in \mathcal{E}$ is an ordered pair of vertices $i \in \mathcal{V}$ and $j \in \mathcal{V}$. It is also denoted $i \rightarrow j$, to indicate that j depends on i . For representing inputs and outputs, it will be convenient to use **incoming half-edges** $\rightarrow j$ and **outgoing half-edges** $i \rightarrow$.

In a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, the **parents** of a vertex j is the set of nodes pointing to j , denoted $\text{pa}(j) := \{i : i \rightarrow j\}$. The **children** of a vertex i is the set of nodes i is pointing to, that is, $\text{ch}(i) := \{j : i \rightarrow j\}$. Vertices without parents are called **roots** and vertices without children are called **leaves**.

A **path** from i to j is defined by a sequence of vertices j_1, \dots, j_m , potentially empty, such that $i \rightarrow j_1 \rightarrow \dots \rightarrow j_m \rightarrow j$. An **acyclic** graph is a graph such that there exists no vertex i with a path from i to i . A **directed acyclic graph** (DAG) is a graph that is both directed and acyclic.

The edges of a DAG define a **partial order** of the vertices, denoted $i \preceq j$ if there exists a path from i to j . The order is partial, since two vertices may not necessarily be linked to each other by a path. Nevertheless, we can define a total order called a **topological order**: any order such that $i \leq j$ if and only if there is no path from j to i .

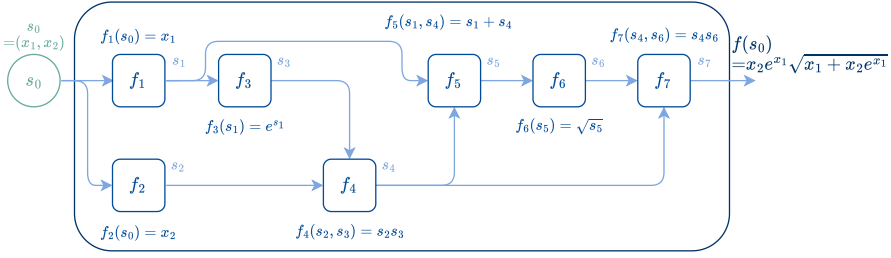


Figure 4.3: Representation of $f(x_1, x_2) = x_2e^{x1}\sqrt{x1+x2}e^{x1}$ as a DAG, with functions and variables as nodes. Edges indicate function and variable dependencies. The function f is decomposed as 8 elementary functions in topological order.

4.1.3 Computer programs as DAGs

We assume that a program defines a mathematically valid function (a.k.a. pure function): the program should return identical values for identical arguments and should not have any side effects. We also assume that the program halts, i.e., that it terminates in a **finite** number of steps. As such a program is made of a finite number of intermediate functions and intermediate variables, the dependencies between functions and variables can be expressed using a directed acyclic graph (DAG). Without loss of generality, we make the following simplifying assumptions:

1. There is a single input $s_0 \in \mathcal{S}_0$.
2. There is a single output $s_K \in \mathcal{S}_K$.
3. Each intermediate function f_k in the program outputs a single variable $s_k \in \mathcal{S}_k$.

We number the nodes as $\mathcal{V} := \{0, 1, \dots, K\}$. Node 0 is the root, corresponding to the input $s_0 \in \mathcal{S}_0$. Node K is the leaf, corresponding to the final output $s_K \in \mathcal{S}_K$. Because of the third assumption above, apart from s_0 , each variable s_k is in **bijection** with a function f_k . Therefore, node 0 represents the input s_0 , and nodes $1, \dots, K$ represent both a function f_k and an output variable s_k .

Edges in the DAG represent dependencies. The parents $i_1, \dots, i_{p_k} := \text{pa}(k)$ of node k , where $p_k := |\text{pa}(k)|$, indicate the variables $s_{\text{pa}(k)} := s_{i_1}, \dots, s_{i_{p_k}}$ that the function f_k needs to perform its computation. Put

Algorithm 4.1 Executing a program**Functions:** f_1, \dots, f_K in topological order**Input:** input $s_0 \in \mathcal{S}_0$ 1: **for** $k := 1, \dots, K$ **do**2: Retrieve parent nodes $(i_1, \dots, i_{p_k}) := \text{pa}(k)$ 3: Compute $s_k := f_k(s_{\text{pa}(k)}) := f_k(s_{i_1}, \dots, s_{i_{p_k}})$ 4: **Output:** $f(s_0) := s_K$

differently, the parents i_1, \dots, i_{p_k} indicate the functions $f_{i_1}, \dots, f_{i_{p_k}}$ that need to be evaluated, prior to evaluating f_k . An example of computation graph in our formalism is presented in Fig. 4.3.

Executing a program

To execute a program, we need to ensure that we evaluate the intermediate functions in the correct order. Therefore, we assume that the nodes $0, 1, \dots, K$ are in a topological order (if this is not the case, we need to perform a topological sort first). We can then execute a program by evaluating for $k \in [K]$

$$s_k := f_k(s_{\text{pa}(k)}) := f_k(s_{i_1}, \dots, s_{i_{p_k}}) \in \mathcal{S}_k.$$

Note that we can either view f_k as a single-input function of $s_{\text{pa}(k)}$, which is a tuple of elements, or as a multi-input function of $s_{i_1}, \dots, s_{i_{p_k}}$. The two views are essentially equivalent.

The procedure for executing a program is summarized in Algorithm 4.1.

Dealing with multiple program inputs or outputs

When a program has multiple inputs, we can always group them into $s_0 \in \mathcal{S}_0$ as $s_0 = (s_{0,1}, \dots, s_{0,N_0})$ with $\mathcal{S}_0 = (\mathcal{S}_{0,1} \times \dots \times \mathcal{S}_{0,N_0})$, since later functions can always filter out what elements of s_0 they need. Likewise, if an intermediate function f_k has multiple outputs, we can always group them as a single output $s_k = (s_{k,1}, \dots, s_{k,N_k})$ with $\mathcal{S}_k = (\mathcal{S}_{k,1} \times \dots \times \mathcal{S}_{k,N_k})$, since later functions can filter out the elements of s_k that they need.



Figure 4.4: Two possible representations of a program. **Left:** Functions and output variables are represented by the same nodes. **Right:** functions and variables are represented by a disjoint set of nodes.

Alternative representation: bipartite graphs

In our formalism, because a function f_k always has a single output s_k , a node k can be seen as representing both the variable s_k and the function f_k . Alternatively, as shown in Fig. 4.4, we can represent variables and functions as separate nodes, that is, using a **bipartite graph**. This formalism is akin to **factor graphs** (Frey *et al.*, 1997; Loeliger, 2004) used in probabilistic modeling, but with directed edges. One advantage of this formalism is that it allows functions to explicitly have multiple outputs. We focus on our formalism for simplicity.

4.1.4 Arithmetic circuits

Arithmetic circuits are one of the simplest examples of computation graph, originating from **computational complexity theory**. Formally, an arithmetic circuit over a field \mathbb{F} , such as the reals \mathbb{R} , is a directed acyclic graph (DAG) whose root nodes are elements of \mathbb{F} and whose functions f_k are either $+$ or \times . The latter are often called **gates**. Contrary to the general computation graph case, because each f_k is either $+$ or \times , it is important to allow the graph to have several root nodes. Root nodes can be either variables or constants, and should belong to \mathbb{F} .

Arithmetic circuits can be used to compute **polynomials**. There are potentially multiple arithmetic circuits for representing a given polynomial. One important question is then to find the most efficient arithmetic circuit for computing a given polynomial. To compare arithmetic circuits representing the same polynomial, an intuitive notion of complexity is the **circuit size**, as defined below.

Definition 4.1 (Circuit and polynomial sizes). The size $S(\mathcal{C})$ of a circuit \mathcal{C} is the number of edges in the directed acyclic graph representing \mathcal{C} . The size $S(f)$ of a polynomial f is the smallest $S(\mathcal{C})$ among all \mathcal{C} representing f .

For more information on arithmetic circuits, we refer the reader to the monograph of Chen *et al.* (2011).

4.2 Feedforward networks

A feedforward network can be seen as a computation chain with **parameterized** functions f_k ,

$$\begin{aligned} s_0 &:= \mathbf{x} \\ s_1 &:= f_1(s_0, \mathbf{w}_1) \\ s_2 &:= f_2(s_1, \mathbf{w}_2) \\ &\vdots \\ s_K &:= f_K(s_{K-1}, \mathbf{w}_K), \end{aligned}$$

for a given input $\mathbf{x} \in \mathcal{X}$ and **learnable parameters** $\mathbf{w}_1, \dots, \mathbf{w}_K \in \mathcal{W}_1 \times \dots \times \mathcal{W}_K$. Each function f_k is called a **layer** and each $s_k \in \mathcal{S}_k$ can be seen as an **intermediate representation** of the input \mathbf{x} . The dimensionality of \mathcal{S}_k is known as the **width** (or number of hidden units) of layer k . A feedforward network defines a function $s_K =: f(\mathbf{x}, \mathbf{w})$ from $\mathcal{X} \times \mathcal{W}$ to \mathcal{S}_K , where $\mathbf{w} := (\mathbf{w}_1, \dots, \mathbf{w}_K) \in \mathcal{W} := \mathcal{W}_1 \times \dots \times \mathcal{W}_K$.

Given such a parameterized program, we can learn the parameters by adjusting \mathbf{w} to fit some data. For instance, given a dataset of $(\mathbf{x}_i, \mathbf{y}_i)$ pairs, we may minimize the squared loss $\|\mathbf{y}_i - f(\mathbf{x}_i, \mathbf{w})\|_2^2$ on average over the data, w.r.t. \mathbf{w} . The minimization of such a loss requires accessing its gradients with respect to \mathbf{w} .

4.3 Multilayer perceptrons

4.3.1 Combining affine layers and activations

In the previous section, we did not specify how to parametrize the feedforward network. A typical parametrization, called the multilayer

perceptron (MLP), uses **fully-connected** (also called **dense**) layers of the form

$$\mathbf{s}_k = f_k(\mathbf{s}_{k-1}, \mathbf{w}_k) := a_k(\mathbf{W}_k \mathbf{s}_{k-1} + \mathbf{b}_k),$$

where we defined the tuple $\mathbf{w}_k := (\mathbf{W}_k, \mathbf{b}_k)$ and where we assumed that \mathbf{W}_k and \mathbf{b}_k are a matrix and vector of appropriate size. We can further decompose the layer into two functions. The function $\mathbf{s} \mapsto \mathbf{W}_k \mathbf{s} + \mathbf{b}_k$ is called an affine layer. The function $\mathbf{v} \mapsto a_k(\mathbf{v})$ is a parameter-free **nonlinearity**, often called an **activation function** (see Section 4.4). The value $\boldsymbol{\alpha}_k := \mathbf{W}_k \mathbf{s}_{k-1} + \mathbf{b}_k$ is often called the **pre-activation value** and the value $a_k(\boldsymbol{\alpha}_k)$ is the **activation value**.

More generally, we may replace the matrix-vector product $\mathbf{W}_k \mathbf{s}_{k-1}$ by any parametrized linear function of \mathbf{s}_{k-1} . For example, **convolutional layers** use the convolution of an input \mathbf{s}_{k-1} with some filters W_k , seen as a linear map.

Remark 4.1 (Dealing with multiple inputs). Sometimes, it is necessary to deal with networks of multiple inputs. For example, suppose we want to design a function $g(\mathbf{x}_1, \mathbf{x}_2, \mathbf{w}_g)$, where $\mathbf{x}_1 \in \mathcal{X}_1$ and $\mathbf{x}_2 \in \mathcal{X}_2$. A simple way to do so is to use the concatenation $\mathbf{x} := \mathbf{x}_1 \oplus \mathbf{x}_2 \in \mathcal{X}_2 \oplus \mathcal{X}_2$ as input to a network $f(\mathbf{x}, \mathbf{w}_f)$. Alternatively, instead of concatenating \mathbf{x}_1 and \mathbf{x}_2 at the input layer, they can be concatenated after having been through one or more hidden layers.

4.3.2 Link with generalized linear models

When the depth is $K = 1$ (only one layer), the output of an MLP is

$$\mathbf{s}_1 = a_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1).$$

This is called a **generalized linear model** (GLM); see Section 3.4. Therefore, MLPs include GLMs as a special case. In particular, when a_1 is the softargmax (see Section 4.4), we obtain (multiclass) logistic regression. For general depth K , the output of an MLP is

$$\mathbf{s}_K = a_K(\mathbf{W}_K \mathbf{s}_{K-1} + \mathbf{b}_K).$$

This can be seen as a GLM on top of **learned representation** \mathbf{s}_{K-1} of the input \mathbf{x} . This is the main appeal of MLPs: they learn the feature

representation and the output model at the same time! We will see that MLPs can also be used as subcomponents in other architectures.

4.4 Activation functions

As we saw in Section 4.3, feedforward networks typically use an activation function a_k at each layer. In this section, we present various nonlinearities from scalar to scalar or from vector to scalar. We also present probability mappings that can be used as such activations.

4.4.1 ReLU and softplus

Many activations are **scalar-to-scalar** functions, but they can also be applied to vectors in an element-wise fashion. The **ReLU** (rectified linear unit) is a popular nonlinearity defined as the **non-negative part** of its input

$$\text{relu}(u) := \max(u, 0) = \begin{cases} u, & u \geq 0 \\ 0, & u < 0 \end{cases}.$$

It is a piecewise linear function and includes a kink at $u = 0$. A multilayer perceptron with ReLU activations is called a **rectifier neural network**. The layers take the form

$$\mathbf{s}_k = \text{relu}(\mathbf{A}_k \mathbf{s}_{k-1} + \mathbf{b}_k),$$

where the ReLU is applied element-wise. The ReLU can be replaced with a smooth approximation (i.e., without kinks), called the **softplus**

$$\text{softplus}(u) := \log(1 + e^u).$$

Unlike the ReLU, it is always strictly positive. Other smoothed variants of the ReLU are possible, see Section 13.4.

4.4.2 Max pooling and log-sum-exp

Many activations are **vector-to-scalar** functions: they **reduce** vectors to a scalar value. This scalar value can be seen as a statistic, “summarizing” the vector.

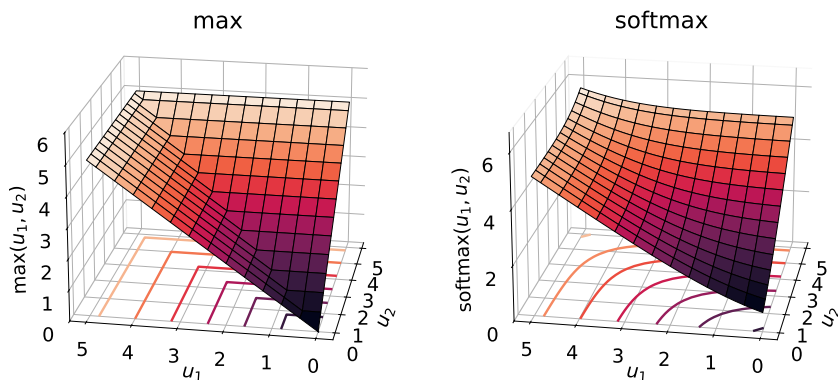


Figure 4.5: The maximum operator is a piecewise linear function. The log-sum-exp (a.k.a. softmax) is a smoothed approximation.

Max pooling

An example of vector-to-scalar reduction is the maximum value, also known as **max pooling**. Given a vector $\mathbf{u} \in \mathbb{R}^M$, it is defined as

$$\max(\mathbf{u}) := \max_{j \in [M]} u_j.$$

Log-sum-exp as a soft maximum

Another example of vector-to-scalar reduction is the **log-sum-exp**,

$$\text{logsumexp}(\mathbf{u}) := \text{softmax}(\mathbf{u}) := \log \sum_{j=1}^M e^{u_j}.$$

As illustrated in Fig. 4.5, it is known to behave like a **soft maximum**. The log-sum-exp can be seen as a generalization of the softplus, as we have for all $u \in \mathbb{R}$

$$\text{logsumexp}((u, 0)) = \text{softplus}(u).$$

A numerically stable implementation of the log-sum-exp is given by

$$\text{logsumexp}(\mathbf{u}) = \text{logsumexp}(\mathbf{u} - c \mathbf{1}) + c,$$

where $c := \max_{j \in [M]} u_j$.

More generally, we can introduce a temperature parameter $\gamma > 0$

$$\text{logsumexp}_\gamma(\mathbf{u}) = \gamma \cdot \text{logsumexp}(\mathbf{u}/\gamma).$$

It can be shown that for all $\mathbf{u} \in \mathbb{R}^M$,

$$\max(\mathbf{u}) \leq \text{logsumexp}_\gamma(\mathbf{u}) \leq \max(\mathbf{u}) + \gamma \cdot \log(M).$$

Therefore, $\text{logsumexp}_\gamma(\mathbf{u}) \rightarrow \max(\mathbf{u})$ as $\gamma \rightarrow 0$. Other definitions of soft maximum are possible; see Section 13.5.

Log-sum-exp as a log-domain sum

Besides its use as a soft maximum, the log-sum-exp often arises for computing sums in the log domain. Indeed, suppose we want to compute $s := \sum_{j=1}^M u_i$, where $u_i > 0$. If we define $\tilde{u}_i := \log u_i$ and $\tilde{s} := \log s$, we then have

$$\tilde{s} = \log \sum_{j=1}^M \exp(\tilde{u}_i).$$

Written differently, we have the identity

$$\log \left(\sum_{i=1}^M u_i \right) = \text{logsumexp}(\log(\mathbf{u})).$$

We can therefore see the log-sum-exp as the sum counterpart of the identity for products

$$\log \left(\prod_{j=1}^M u_i \right) = \sum_{i=1}^M \log(u_i).$$

As an example, we use the log-sum-exp to perform the forward-backward algorithm in the log-domain in Section 10.7.1.

4.4.3 Sigmoids: binary step and logistic functions

Oftentimes, we want to map a real value to a number in $[0, 1]$, that can represent the probability of an event. For that purpose, we generally use **sigmoids**. A sigmoid is a function with a characteristic “S”-shaped curve. These functions are **scalar-to-scalar** probability mappings: they are used to squash real values to $[0, 1]$.

Binary step function

An example is the **binary step function**, also known as **Heaviside step function**,

$$\text{step}(u) := \begin{cases} 1, & u \geq 0 \\ 0, & u < 0 \end{cases}.$$

It is a mapping from \mathbb{R} to $\{0, 1\}$. Unfortunately, it has a discontinuity: a jump in its graph at $u = 0$. Moreover, because the function is constant at all other points, it has zero derivative at these points, which makes it difficult to use as part of a neural network trained with backpropagation.

Logistic function

A better sigmoid is the **logistic function**, which is a mapping from \mathbb{R} to $(0, 1)$ and is defined as

$$\begin{aligned} \text{logistic}(u) &:= \frac{1}{1 + e^{-u}} \\ &= \frac{e^u}{1 + e^u} \\ &= \frac{1}{2} + \frac{1}{2} \tanh\left(\frac{u}{2}\right). \end{aligned}$$

It maps $(-\infty, 0)$ to $(0, 0.5)$, $[0, +\infty)$ to $[0.5, 1)$ and it satisfies $\text{logistic}(0) = 0.5$. It can therefore be seen as mapping from real values to probability values. The logistic can be seen as a differentiable approximation to the discontinuous binary step function $\text{step}(u)$. The logistic function can be shown to be the derivative of *softplus*, i.e., for all $u \in \mathbb{R}$

$$\text{softplus}'(u) = \text{logistic}(u).$$

Two important properties of the logistic function are that for all $u \in \mathbb{R}$

$$\text{logistic}(-u) = 1 - \text{logistic}(u)$$

and

$$\begin{aligned} \text{logistic}'(u) &= \text{logistic}(u) \cdot \text{logistic}(-u) \\ &= \text{logistic}(u) \cdot (1 - \text{logistic}(u)). \end{aligned}$$

Other sigmoids are possible; see Section [13.6](#).

4.4.4 Probability mappings: argmax and softargmax

It is often useful to transform a real vector into a vector of probabilities. This is a mapping from \mathbb{R}^M to the probability simplex, defined by

$$\Delta^M := \left\{ \boldsymbol{\pi} \in \mathbb{R}^M : \forall j \in [M], \pi_j \geq 0, \sum_{j=1}^M \pi_j = 1 \right\}.$$

Two examples of such **vector-to-vector** probability mappings are the argmax and the softargmax.

Argmax

The argmax operator is defined by

$$\text{argmax}(\mathbf{u}) := \phi \left(\arg \max_{j \in [M]} u_j \right) \in \{\mathbf{e}_1, \dots, \mathbf{e}_M\},$$

where $\phi(j)$ denotes the one-hot encoding of an integer $j \in [M]$, that is,

$$\phi(j) := (0, \dots, 0, \underbrace{1}_j, 0, \dots, 0) = \mathbf{e}_j \in \{0, 1\}^M.$$

This mapping puts all the probability mass onto a single coordinate (in case of ties, we pick a single coordinate arbitrarily). Unfortunately, this mapping is a discontinuous function.

Softargmax

As a differentiable everywhere relaxation, we can use the **softargmax** defined by

$$\text{softargmax}(\mathbf{u}) := \frac{\exp(\mathbf{u})}{\sum_{j=1}^M \exp(u_j)} \in \text{relint}(\Delta^M).$$

This operator is commonly known in the literature as *softmax* but this is a misnomer: this operator really defines a differentiable relaxation of the argmax. The output of the softargmax belongs to the relative interior of the probability simplex $\text{relint}(\Delta^M) = \{\boldsymbol{\pi} \in \Delta^M : \boldsymbol{\pi} > \mathbf{0}\}$, meaning that it can never reach the borders of the simplex. If we denote

$\boldsymbol{\pi} = \text{softargmax}(\mathbf{u})$, this means that $\pi_j \in (0, 1)$, that is, π_j can never be exactly 0 or 1. The softargmax is the gradient of log-sum-exp,

$$\nabla \text{logsumexp}(\mathbf{u}) = \text{softargmax}(\mathbf{u}).$$

The softargmax can be seen as a generalization of the logistic function, as we have for all $u \in \mathbb{R}$

$$[\text{softargmax}((u, 0))]_1 = \text{logistic}(u).$$

Remark 4.2 (Degrees of freedom and invertibility of softargmax). The softargmax operator satisfies the property for all $\mathbf{u} \in \mathbb{R}^M$ and $c \in \mathbb{R}$

$$\boldsymbol{\pi} := \text{softargmax}(\mathbf{u}) = \text{softargmax}(\mathbf{u} + c \mathbf{1}).$$

This means that the softargmax operator has $M - 1$ degrees of freedom and is a non-invertible function. However, due to the above property, without loss of generality, we can impose $\mathbf{u}^\top \mathbf{1} = \sum_{i=1}^M u_i = 0$ (if this is not the case, we simply do $u_i \leftarrow u_i - \bar{u}$, where $\bar{u} := \frac{1}{M} \sum_{j=1}^M u_j$). Using this constraint together with

$$\log \pi_i = u_i - \log \sum_{j=1}^M \exp(u_j),$$

we then obtain

$$\sum_{i=1}^M \log \pi_i = -M \log \sum_{j=1}^M \exp(u_j)$$

so that

$$u_i = [\text{softargmax}^{-1}(\boldsymbol{\pi})]_i = \log \pi_i - \frac{1}{M} \sum_{j=1}^M \log \pi_j.$$

4.5 Normalization layers

Intermediate states in neural networks, such as activations, can often attain a wide range of different values, potentially making it difficult for gradient descent to converge. To remedy this issue, we can introduce

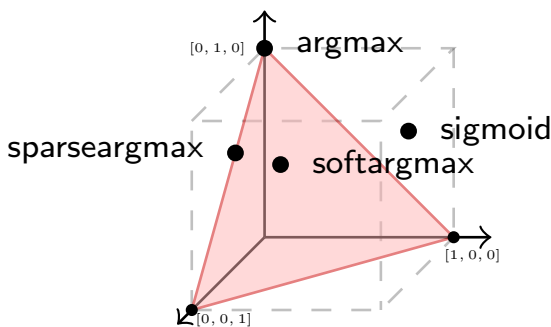


Figure 4.6: The argmax, the softargmax and the sparseargmax (see Section 13.7) are mapping from \mathbb{R}^M to the probability simplex $\Delta^M := \{\pi \in \mathbb{R}_+^M : \langle \pi, \mathbf{1} \rangle = 1\}$. More precisely, the argmax maps to vertices of the probability simplex, the set of one-hot vectors $\{e_1, \dots, e_M\}$, the softargmax maps to the relative interior $\text{relint}(\Delta^M) = \{\pi \in \Delta^M : \pi > \mathbf{0}\}$ and the sparseargmax maps to the entire probability simplex Δ^M , including sparse probability vectors. Sigmoids applied element-wise return vectors in the hypercube $[0, 1]^M$, whose vertices are the 2^M vectors in $\{0, 1\}^M$.

normalization layers at suitable locations in the network. In this section, we present the two most popular ones: batch normalization and layer normalization.

4.5.1 Batch normalization

Suppose we are given a batch of intermediate variables (such as activations or states) $\mathbf{s}_1, \dots, \mathbf{s}_B \in \mathbb{R}^D$, where $\mathbf{s}_i := (s_{i,1}, \dots, s_{i,D})$, obtained by applying some function to B samples drawn from the training set (we omit the dependency on the layer index k for clarity). In batch normalization (Ioffe and Szegedy, 2015), we normalize the values by calculating the **standard score** (a.k.a. **z-score**) across batch samples,

$$\begin{aligned} \mu_j &:= \frac{1}{B} \sum_{i=1}^B s_{i,j} \quad \forall j \in [D] \\ \sigma_j^2 &:= \frac{1}{B} \sum_{i=1}^B (s_{i,j} - \mu_j)^2 \quad \forall j \in [D] \\ \hat{s}_{i,j} &:= \frac{s_{i,j} - \mu_j}{\sigma_j} \quad i \in [B], j \in [D]. \end{aligned}$$

Here, the means μ_j and the standard deviations σ_j are computed for each feature $j \in [D]$ across samples $i \in [B]$ in the batch. In practice, we often add a small value $\varepsilon > 0$ to σ_j to avoid division by zero or numerical instabilities. Moreover, we often rescale the values as

$$\tilde{s}_{i,j} := \beta_j + \gamma_j \hat{s}_{i,j} \quad i \in [B], j \in [D],$$

where the means $\beta := (\beta_1, \dots, \beta_D)$ and standard deviations $\gamma := (\gamma_1, \dots, \gamma_D)$ are learnable parameters.

One issue with batch normalization is that the means μ_j and standard deviations σ_j cannot be computed at inference time, as there is no notion of training batch (a single sample would lead to a variance of 0). To address this issue, we can estimate means $\hat{\mu}_j$ and standard deviations $\hat{\sigma}_j$ across the whole training set during the course of training, usually using a **running average**. A practical batch normalization implementation therefore needs to maintain D mean and standard deviation statistics, so as to be able to use them at inference time.

4.5.2 Layer normalization

As an alternative, in layer normalization (Ba *et al.*, 2016), we instead standardize the values by summing across features,

$$\begin{aligned} \mu_i &:= \frac{1}{D} \sum_{j=1}^D s_{i,j} \quad \forall i \in [B] \\ \sigma_i^2 &:= \frac{1}{D} \sum_{j=1}^D (s_{i,j} - \mu_i)^2 \quad \forall i \in [B] \\ \hat{s}_{i,j} &:= \frac{s_{i,j} - \mu_i}{\sigma_i} \quad i \in [B], j \in [D]. \end{aligned}$$

This time, the means μ_i and the standard deviations σ_i are computed for each sample $i \in [B]$ across features $j \in [D]$ (as before, we often add a small value ε to σ_i). Similarly to batch normalization, we often rescale the values as

$$\tilde{s}_{i,j} := \beta_j + \gamma_j \hat{s}_{i,j} \quad i \in [B], j \in [D],$$

where the means $\beta := (\beta_1, \dots, \beta_D)$ and standard deviations $\gamma := (\gamma_1, \dots, \gamma_D)$ are learnable parameters. A key advantage of layer normalization compared to batch normalization is that it is well defined

at inference time, since it is applied on a per-sample basis and does not rely on the notion of training batch. As a result, we can view layer normalization as a function that to any $\mathbf{s}_i \in \mathbb{R}^D$ (regardless of whether it is part of a batch or not) associates

$$\tilde{\mathbf{s}}_i := \text{LayerNorm}(\mathbf{s}_i).$$

4.6 Residual neural networks

We now discuss another feedforward network parametrization: residual neural networks. Consider a feedforward network with $K + 1$ layers f_1, \dots, f_K, f_{K+1} . Surely, as long as f_{K+1} can exactly represent the identity function, the set of functions that this feedforward network can express should be a superset of the functions that f_1, \dots, f_K can express. In other words, depth should in theory not hurt the expressive power of feedforward networks. Unfortunately, the assumption that each f_k can exactly represent the identity function may not hold in practice. This means that deeper networks can sometimes be more difficult to train than shallower ones, making the accuracy saturate or degrade as a function of depth.

The key idea of residual neural networks (He *et al.*, 2016) is to design layers f_k , called **residual blocks**, that make it easier to represent the identity function. Formally, a residual block takes the form

$$\mathbf{s}_k = f_k(\mathbf{s}_{k-1}, \mathbf{w}_k) := \mathbf{s}_{k-1} + h_k(\mathbf{s}_{k-1}, \mathbf{w}_k).$$

The function h_k is called **residual**, since it models the difference $\mathbf{s}_k - \mathbf{s}_{k-1}$. The addition with \mathbf{s}_{k-1} is often called a **skip connection**. As long as it is easy to adjust \mathbf{w}_k so that $h_k(\mathbf{s}_{k-1}, \mathbf{w}_k) = \mathbf{0}$, f_k can freely become the identity function. For instance, if we use

$$h_k(\mathbf{s}_{k-1}, \mathbf{w}_k) := \mathbf{C}_k a_k(\mathbf{W}_k \mathbf{s}_{k-1} + \mathbf{b}_k) + \mathbf{d}_k,$$

where $\mathbf{w}_k := (\mathbf{W}_k, \mathbf{b}_k, \mathbf{C}_k, \mathbf{d}_k)$, it suffices to set \mathbf{C}_k and \mathbf{d}_k to a zero matrix and vector. Residual blocks are known to remedy the so-called vanishing gradient problem.

Many papers and software packages include an additional activation and instead define the residual block as

$$\mathbf{s}_k = f_k(\mathbf{s}_{k-1}, \mathbf{w}_k) := a_k(\mathbf{s}_{k-1} + h_k(\mathbf{s}_{k-1}, \mathbf{w}_k)),$$

where a_k is typically chosen to be the ReLU activation. Whether to include this additional activation or not is essentially a modelling choice. In practice, residual blocks may also include additional operations such as batch norm and convolutional layers.

4.7 Recurrent neural networks

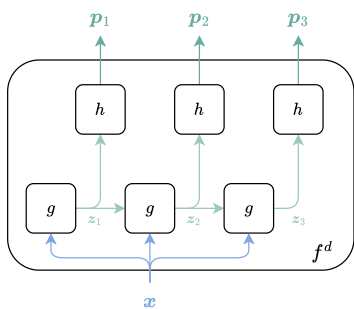
Recurrent neural networks (RNNs) are a class of neural networks that operate on sequences of vectors, either as input, output or both. Their actual parametrization depends on the setup but the core idea is to maintain a **state vector** that is updated from step to step by a recursive function that uses **shared parameters** across steps. Unrolling this recursion defines a valid computational graph, as we will see in Chapter 8. We distinguish between the following setups illustrated in Fig. 4.7:

- Vector to sequence (one to many):
 $f^d: \mathbb{R}^D \times \mathbb{R}^P \rightarrow \mathbb{R}^{L \times M}$
- Sequence to vector (many to one):
 $f^e: \mathbb{R}^{L \times D} \times \mathbb{R}^P \rightarrow \mathbb{R}^M$
- Sequence to sequence (many to many, aligned):
 $f^a: \mathbb{R}^{L \times D} \times \mathbb{R}^P \rightarrow \mathbb{R}^{L \times M}$
- Sequence to sequence (many to many, unaligned):
 $f^u: \mathbb{R}^{L \times D} \times \mathbb{R}^P \rightarrow \mathbb{R}^{L' \times M}$

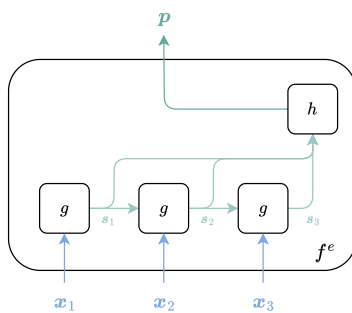
where L stands for length. Note that we use the same number of parameters P for each setup for notational convenience, but this of course does not need to be the case. Throughout this section, we use the notation $\mathbf{p}_{1:L} := (\mathbf{p}_1, \dots, \mathbf{p}_L)$ for a sequence of L vectors.

4.7.1 Vector to sequence

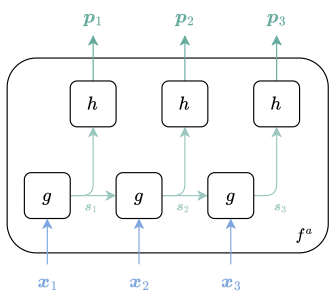
In this setting, we define a **decoder** function $\mathbf{p}_{1:L} = f^d(\mathbf{x}, \mathbf{w})$ from an **input vector** $\mathbf{x} \in \mathbb{R}^D$ and parameters $\mathbf{w} \in \mathbb{R}^P$ to an **output sequence** $\mathbf{p}_{1:L} \in \mathbb{R}^{L \times M}$. This is for instance useful for image caption generation, where a sentence (a sequence of word embeddings) is generated from



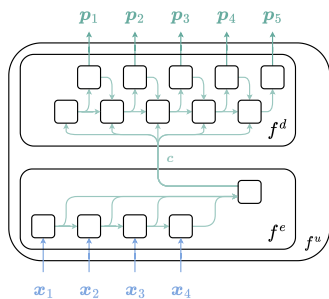
(a) One to many (decoder)



(b) Many to one (encoder)



(c) Sequence to sequence aligned



(d) Sequence to sequence unaligned

Figure 4.7: Recurrent neural network architectures

an image (a vector of pixels). Formally, we may define $\mathbf{p}_{1:L} := f^d(\mathbf{x}, \mathbf{w})$ through the recursion

$$\begin{aligned} \mathbf{z}_l &:= g(\mathbf{x}, \mathbf{z}_{l-1}, \mathbf{w}_g) \quad l \in [L] \\ \mathbf{p}_l &:= h(\mathbf{z}_l, \mathbf{w}_h) \quad l \in [L]. \end{aligned}$$

where $\mathbf{w} := (\mathbf{w}_g, \mathbf{w}_h, \mathbf{z}_0)$. The goal of g is to update the current **decoder state** \mathbf{z}_l given the input \mathbf{x} , and the previous decoder state \mathbf{z}_{l-1} . The goal of h is to generate the output \mathbf{p}_l given the current decoder state \mathbf{z}_l . Importantly, the parameters of g and h are **shared across steps**. Typically, g and h are parametrized using one-hidden-layer MLPs. Note that g has multiple inputs; we discuss how to deal with such cases in Section 4.3.

4.7.2 Sequence to vector

In this setting, we define an **encoder** function $\mathbf{p} = f^e(\mathbf{x}_{1:L}, \mathbf{w})$ from an **input sequence** $\mathbf{x}_{1:L} \in \mathbb{R}^{L \times D}$ and parameters $\mathbf{w} \in \mathbb{R}^P$ to an **output vector** $\mathbf{p} \in \mathbb{R}^M$. This is for instance useful for sequence classification, such as sentiment analysis. Formally, we may define $\mathbf{p} := f^e(\mathbf{x}_{1:L}, \mathbf{w})$ using the recursion

$$\begin{aligned} \mathbf{s}_l &:= \gamma(\mathbf{x}_l, \mathbf{s}_{l-1}, \mathbf{w}_g) \quad l \in [L] \\ \mathbf{p} &= \text{pooling}(\mathbf{s}_{1:L}) \end{aligned}$$

where $\mathbf{w} := (\mathbf{w}_g, \mathbf{s}_0)$. The goal of γ is similar as g , except that it updates **encoder states** and does not take previous predictions as input. The **pooling** function is typically parameter-less. Its goal is to reduce a sequence to a vector. Examples include using the last state, the average of states and the coordinate-wise maximum of states.

4.7.3 Sequence to sequence (aligned)

In this setting, we define a function $\mathbf{p}_{1:L} = f^a(\mathbf{x}_{1:L}, \mathbf{w})$ from an **input sequence** $\mathbf{x}_{1:L} \in \mathbb{R}^{L \times D}$ and parameters $\mathbf{w} \in \mathbb{R}^P$ to an **output sequence** $\mathbf{p}_{1:L} \in \mathbb{R}^{L \times M}$, which we assume to be of the **same length**. An example of application is part-of-speech tagging, where the goal is to assign each word \mathbf{x}_l to a part-of-speech (noun, verb, adjective, etc).

Formally, we may define $\mathbf{p}_{1:L} = f^a(\mathbf{x}_{1:L}, \mathbf{w})$ as

$$\begin{aligned} \mathbf{s}_l &:= \gamma(\mathbf{x}_l, \mathbf{s}_{l-1}, \mathbf{w}_\gamma) \quad l \in [L] \\ \mathbf{p}_l &:= h(\mathbf{s}_l, \mathbf{w}_h) \quad l \in [L] \end{aligned}$$

where $\mathbf{w} := (\mathbf{w}_\gamma, \mathbf{w}_h, \mathbf{s}_0)$. The function γ and h are similar as before.

4.7.4 Sequence to sequence (unaligned)

In this setting, we define a function $\mathbf{p}_{1:L'} = f^u(\mathbf{x}_{1:L}, \mathbf{w})$ from an **input sequence** $\mathbf{x}_{1:L} \in \mathbb{R}^{L \times D}$ and parameters $\mathbf{w} \in \mathbb{R}^P$ to an output sequence $\mathbf{p}_{1:L'} \in \mathbb{R}^{L' \times M}$, which potentially has a **different length**. An example of application is machine translation, where the sentences in the source and target languages do not necessarily have the same length. Typically, $\mathbf{p}_{1:L'} = f^u(\mathbf{x}_{1:L}, \mathbf{w})$ is defined as the following two steps

$$\begin{aligned} \mathbf{c} &:= f^e(\mathbf{x}_{1:L}, \mathbf{w}_e) \\ \mathbf{p}_{1:L'} &:= f^d(\mathbf{c}, \mathbf{w}_d) \end{aligned}$$

where $\mathbf{w} := (\mathbf{w}_e, \mathbf{w}_d)$, and where we reused the previously-defined encoder f_e and decoder f_d . Putting the two steps together, we obtain

$$\begin{aligned} \mathbf{s}_l &:= \gamma(\mathbf{x}_l, \mathbf{s}_{l-1}, \mathbf{w}_\gamma) \quad l \in [L] \\ \mathbf{c} &= \text{pooling}(\mathbf{s}_{1:L}) \\ \mathbf{z}_l &:= g(\mathbf{c}, \mathbf{p}_{l-1}, \mathbf{z}_{l-1}, \mathbf{w}_g) \quad l \in [L'] \\ \mathbf{p}_l &:= h(\mathbf{z}_l, \mathbf{w}_h) \quad l \in [L']. \end{aligned}$$

This architecture is aptly named the **encoder-decoder** architecture. Note that we denoted the length of the target sequence as L' . However, in practice, the target length can be input dependent and is often not known ahead of time. To deal with this issue, the vocabulary (of size D is our notation) is typically augmented with an “end of sequence” (EOS) token so that, at inference time, we know when to stop generating the output sequence. One disadvantage of this encoder-decoder architecture, however, is that all the information about the input sequence is contained in the **context vector** \mathbf{c} , which can therefore become a **bottleneck**.

4.8 Transformers

Transformers (Vaswani *et al.*, 2017) are one of the most successful recent developments in deep learning. In this section, we review Transformers, component by component.

4.8.1 Attention

The goal of an attention layer is to map a sequence of inputs $\mathbf{v}_1, \dots, \mathbf{v}_L \in \mathbb{R}^{D_v}$ to a sequence of outputs $\mathbf{u}_1, \dots, \mathbf{u}_L \in \mathbb{R}^{D_v}$, of same dimension. A natural idea is to define each output element \mathbf{u}_i using a linear combination of the input elements,

$$\mathbf{u}_i := \sum_{j=1}^L a_{i,j} \mathbf{v}_j \in \mathbb{R}^{D_v}.$$

We typically use a convex combination: we assume that the combination weights are such that $\mathbf{a}_i := (a_{i,1}, \dots, a_{i,L}) \in \Delta^L$. This ensures that increasing a coefficient $a_{i,j}$ is made at the expense of decreasing another coefficient $a_{i,j'}$, for $j \neq j'$. Let us form the matrices $\mathbf{V} \in \mathbb{R}^{L \times D_v}$ and $\mathbf{U} \in \mathbb{R}^{L \times D_v}$ by stacking $\mathbf{v}_1, \dots, \mathbf{v}_L$ and $\mathbf{u}_1, \dots, \mathbf{u}_L$, seen as row vectors. Let us also form the **attention matrix** $\mathbf{A} \in \Delta^{L \times L}$ gathering the entries $a_{i,j}$, where $\Delta^{L \times L}$ denotes the set of row-wise stochastic $L \times L$ matrices. Then, we can rewrite attention succinctly as

$$\mathbf{U} = \mathbf{A}\mathbf{V} \in \mathbb{R}^{L \times D_v}.$$

Following Section 6.2, we can view attention as a **soft dictionary lookup**. From this perspective, the matrix $\mathbf{V} \in \mathbb{R}^{L \times D_v}$ plays the role of dictionary values, and the attention matrix can be defined as

$$\mathbf{A} := \text{softargmax}(\mathbf{Q}\mathbf{K}^\top) \in \Delta^{L \times L},$$

where $\mathbf{Q} \in \mathbb{R}^{L \times D_k}$ and $\mathbf{K} \in \mathbb{R}^{L \times D_k}$ play the roles of queries and dictionary keys, respectively. Intuitively, $\mathbf{Q}\mathbf{K}^\top \in \mathbb{R}^{L \times L}$ can be seen as a similarity matrix containing the similarities between queries and dictionary keys. Putting everything together, we can define attention as

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) := \text{softargmax}(\mathbf{Q}\mathbf{K}^\top) \mathbf{V} \in \mathbb{R}^{L \times D_v},$$

In **masked attention**, we additionally incorporate a mask $\mathbf{M} \in \mathbb{R}^{L \times L}$,

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}; \mathbf{M}) := \text{softargmax} \left((\mathbf{Q}\mathbf{K}^\top) \circ \mathbf{M} \right) \mathbf{V} \in \mathbb{R}^{L \times D_v},$$

where \circ denotes the Hadamard product (element-wise multiplication). The mask can be used to force some attention weights $a_{i,j}$ to be zero, by setting the corresponding mask entry $m_{i,j}$ to $-\infty$. For instance, in decoder-only architectures (Section 4.8.8), the mask will prove useful to define **causal attention** for autoregressive models.

Remark 4.3 (Scaled attention). Practical implementations often use **scaled attention**, where a factor of $\frac{1}{\sqrt{D_k}}$ is used within the softargmax, in order to reduce the variance (Vaswani *et al.*, 2017, Footnote 4). We omit this detail for clarity.

4.8.2 Self-attention

Suppose we are given a sequence of feature vectors $\mathbf{x}_1, \dots, \mathbf{x}_L \in \mathbb{R}^D$, that we gather as a matrix $\mathbf{X} \in \mathbb{R}^{L \times D}$. To define the attention weights $a_{i,j}$, a natural idea is to consider the similarity between \mathbf{x}_i and \mathbf{x}_j , as measured by the inner product $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$. To ensure that $\mathbf{a}_i := (a_{i,1}, \dots, a_{i,L}) \in \Delta^L$, we can then define

$$a_{i,j} := \frac{\exp(\langle \mathbf{x}_i, \mathbf{x}_j \rangle)}{\sum_{j'=1}^L \exp(\langle \mathbf{x}_i, \mathbf{x}_{j'} \rangle)} \in (0, 1).$$

In matrix notation, this can be written

$$\mathbf{A} := \text{softargmax}(\mathbf{X}\mathbf{X}^\top) \in \Delta^{L \times L},$$

where softargmax is applied in a row-wise fashion. The matrix $\mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{L \times L}$ is the **Gram matrix** associated with the row vectors $\mathbf{x}_1, \dots, \mathbf{x}_L$. In the notation of the previous section, this corresponds to using $\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$ with $\mathbf{Q} = \mathbf{K} = \mathbf{V} = \mathbf{X}$. In other words, the elements of the sequence “pay attention” to each other. This is called **self-attention**.

So far, the formulation we described is parameter-free. To give more expressive power to the self-attention layer, Vaswani *et al.* (2017)

proposed instead to define \mathbf{Q} , \mathbf{K} and \mathbf{V} by projecting \mathbf{X} as

$$\begin{aligned}\mathbf{Q} &:= \mathbf{X}\mathbf{W}^Q \in \mathbb{R}^{L \times D_k} \\ \mathbf{K} &:= \mathbf{X}\mathbf{W}^K \in \mathbb{R}^{L \times D_k} \\ \mathbf{V} &:= \mathbf{X}\mathbf{W}^V \in \mathbb{R}^{L \times D_v},\end{aligned}$$

using the learned weight matrices

$$\begin{aligned}\mathbf{W}^Q &\in \mathbb{R}^{D \times D_k} \\ \mathbf{W}^K &\in \mathbb{R}^{D \times D_k} \\ \mathbf{W}^V &\in \mathbb{R}^{D \times D_v}.\end{aligned}$$

Importantly, the size of the weight matrices is independent of the length L of the sequences. This allows self-attention to work with sequence of arbitrary length.

4.8.3 Multi-head attention

In order to be able capture multiple patterns of attention, Vaswani *et al.* (2017) found it beneficial to define H attention heads

$$\mathbf{Y}_i := \text{Attention}(\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i; \mathbf{M}) \in \mathbb{R}^{L \times D_v},$$

where

$$\begin{aligned}\mathbf{Q}_i &:= \mathbf{X}\mathbf{W}_i^Q \in \mathbb{R}^{L \times D_k} \\ \mathbf{K}_i &:= \mathbf{X}\mathbf{W}_i^K \in \mathbb{R}^{L \times D_k} \\ \mathbf{V}_i &:= \mathbf{X}\mathbf{W}_i^V \in \mathbb{R}^{L \times D_v},\end{aligned}$$

using the learned weight matrices

$$\begin{aligned}\mathbf{W}_i^Q &\in \mathbb{R}^{D \times D_k} \\ \mathbf{W}_i^K &\in \mathbb{R}^{D \times D_k} \\ \mathbf{W}_i^V &\in \mathbb{R}^{D \times D_v}.\end{aligned}$$

Let us denote the concatenation of the H attention heads by

$$\text{Concat}(\mathbf{Y}_1, \dots, \mathbf{Y}_H) \in \mathbb{R}^{L \times H D_v}.$$

Algorithm 4.2 Multi-head attention with H attention heads

Input: $\mathbf{X} \in \mathbb{R}^{L \times D}$, optional mask $\mathbf{M} \in \mathbb{R}^{L \times L}$
Parameters: $\{\mathbf{W}_i^Q\}_{i=1}^H$, $\{\mathbf{W}_i^K\}_{i=1}^H$, $\{\mathbf{W}_i^V\}_{i=1}^H$, \mathbf{W}^O
 1: **for** $i := 1, \dots, H$ **do**
 2: $\mathbf{Q}_i := \mathbf{X} \mathbf{W}_i^Q \in \mathbb{R}^{L \times D_k}$
 3: $\mathbf{K}_i := \mathbf{X} \mathbf{W}_i^K \in \mathbb{R}^{L \times D_k}$
 4: $\mathbf{V}_i := \mathbf{X} \mathbf{W}_i^V \in \mathbb{R}^{L \times D_v}$
 5: $\mathbf{Y}_i := \text{Attention}(\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i; \mathbf{M}) \in \mathbb{R}^{L \times D_v}$
 6: $\mathbf{Y} := \text{Concat}(\mathbf{Y}_1, \dots, \mathbf{Y}_H) \mathbf{W}^O$
Output: $\text{MultiheadAttention}(\mathbf{X}; \mathbf{M}) := \mathbf{Y} \in \mathbb{R}^{L \times D}$

Vaswani *et al.* (2017) then define multi-head attention as

$$\begin{aligned}
 \text{MultiheadAttention}(\mathbf{X}; \mathbf{M}) &:= \text{Concat}(\mathbf{Y}_1, \dots, \mathbf{Y}_H) \mathbf{W}^O \\
 &= \sum_{i=1}^H \text{Attention}(\mathbf{X} \mathbf{W}_i^Q, \mathbf{X} \mathbf{W}_i^K, \mathbf{X} \mathbf{W}_i^V) \mathbf{W}_i^O \\
 &\in \mathbb{R}^{L \times D}
 \end{aligned}$$

where $\mathbf{W}^O \in \mathbb{R}^{HD_v \times D}$ is a learned matrix. We summarize the procedure in Algorithm 4.2. We use a for loop for the sake of clarity; a GPU-friendly implementation would instead be based on a single matrix multiplication.

4.8.4 Transformer layer

To improve training efficiency, we can introduce residual connections and a first layer normalization (Section 4.5.2), to define

$$\mathbf{Z} := \text{LayerNorm}_1(\text{MultiheadAttention}(\mathbf{X}; \mathbf{M}) + \mathbf{X}) \in \mathbb{R}^{L \times D}.$$

To improve expressiveness, a Transformer layer further uses an MLP and a second layer normalization,

$$\mathbf{X} \leftarrow \text{LayerNorm}_2(\text{MLP}(\mathbf{Z}) + \mathbf{Z}) \in \mathbb{R}^{L \times D}.$$

The subscripts in `LayerNorm` are used to emphasize that the two layers use their own parameters. In addition, normalization is applied in an

Algorithm 4.3 Transformer's MLP block**Input:** $\mathbf{Z} \in \mathbb{R}^{L \times D}$ **Parameters:** $\mathbf{W}_1 \in \mathbb{R}^{D \times D_1}, \mathbf{W}_2 \in \mathbb{R}^{D_1 \times D}$ 1: $\mathbf{Z} \leftarrow \mathbf{Z}\mathbf{W}_1 \in \mathbb{R}^{L \times D_1}$ 2: $\mathbf{Z} \leftarrow \sigma(\mathbf{Z}) \in \mathbb{R}^{L \times D_1}$ 3: $\mathbf{Z} \leftarrow \mathbf{Z}\mathbf{W}_2 \in \mathbb{R}^{L \times D}$ **Output:** $\mathbf{Z} \in \mathbb{R}^{L \times D}$

element-wise (token-wise) fashion. The MLP block typically uses a single hidden layer with an activation function σ , such as a GELU (Gaussian error linear unit); see Algorithm 4.3. Because of the use of residual connections, the input and output dimensions of the MLP block must be the same.

Remark 4.4 (Post-normalization vs. pre-normalization). Our description of the Transformer layer follows the original formulation (Vaswani *et al.*, 2017), which uses post-normalization. Some implementations instead rely on pre-normalization, that is,

$$\begin{aligned} \mathbf{Z} &:= \text{MultiheadAttention}(\text{LayerNorm}_1(\mathbf{X}); \mathbf{M}) + \mathbf{X} \in \mathbb{R}^{L \times D} \\ \mathbf{X} &\leftarrow \text{MLP}(\text{LayerNorm}_2(\mathbf{Z})) + \mathbf{Z} \in \mathbb{R}^{L \times D}. \end{aligned}$$

4.8.5 Transformer block

A Transformer layer can be iterated K times (each time with different parameters) to define a Transformer block of depth K . Keeping the dependency on parameters implicit, we can see a Transformer, with an optional mask $\mathbf{M} \in \mathbb{R}^{L \times L}$, as a function

$$\mathbf{X} \mapsto \text{Transformer}(\mathbf{X}; \mathbf{M}) = \text{Transformer}(\mathbf{x}_1, \dots, \mathbf{x}_L; \mathbf{M})$$

from $\mathbb{R}^{L \times D}$ to $\mathbb{R}^{L \times D}$. The Transformer takes a sequence $\mathbf{X} \in \mathbb{R}^{L \times D}$ and uses the inter-dependencies between sequence elements to produce a representation of that sequence. We summarize the procedure in Algorithm 4.4. The MultiheadAttention, LayerNorm and MLP blocks are indexed by k to emphasize that their parameters are different for each iteration.

Algorithm 4.4 Transformer block of depth K

Input: $\mathbf{X} \in \mathbb{R}^{L \times D}$ optional mask $\mathbf{M} \in \mathbb{R}^{L \times L}$ **Parameters:** Multi-head attention parameters, MLP parameters and layer norm parameters (each depth k uses different parameters)1: **for** $k := 1, \dots, K$ **do**2: $\mathbf{Y} := \text{MultiheadAttention}_k(\mathbf{X}; \mathbf{M}) \in \mathbb{R}^{L \times D}$ 3: $\mathbf{Z} := \text{LayerNorm}_{k,1}(\mathbf{Y} + \mathbf{X}) \in \mathbb{R}^{L \times D}$ \triangleright Residual connection4: $\mathbf{X} \leftarrow \text{LayerNorm}_{k,2}(\text{MLP}_k(\mathbf{Z}) + \mathbf{Z}) \in \mathbb{R}^{L \times D}$ \triangleright MLP layer**Output:** $\mathbf{X} \in \mathbb{R}^{L \times D}$

Number of parameters and computational complexity

The Transformer layer can be seen as a function from $\mathbb{R}^{L \times D}$ to $\mathbb{R}^{L \times D}$. To offer the same function signature, a standard multi-layer perceptron would have needed $O(N^2 D^2)$ parameters and a forward pass through the network would have had a time complexity of $O(N^2 D^2)$ as well. In contrast, a Transformer layer has $O(D^2)$ parameters. Self-attention has a complexity of $O(N^2 D)$ and the final MLP layer has a complexity of $O(ND^2)$.

4.8.6 Token encoding

Text can be represented as a sequence x_1, \dots, x_L of **discrete** symbols, often called tokens. Here, each token $x_i \in [M]$, where M is the vocabulary size, can correspond to words, subwords, or even individual characters, depending on the tokenization procedure. To obtain a sequence of **continuous** vectors $\mathbf{x}_1, \dots, \mathbf{x}_L$, where $\mathbf{x}_i \in \mathbb{R}^D$, it is common to use an **embedding layer**. This layer transforms $x_i \in [M]$ into $\mathbf{x}_i \in \mathbb{R}^D$ using the linear projection

$$\mathbf{x}_i := \mathbf{W}^E \mathbf{e}_{x_i},$$

where $\mathbf{e}_j \in \mathbb{R}^M$ is the one-hot encoding of $j \in [M]$ and $\mathbf{W}^E \in \mathbb{R}^{D \times M}$ is a learnable embedding matrix. The column $\mathbf{W}_{:,j}^E \in \mathbb{R}^D$ can be interpreted as the continuous representation of token $j \in [V]$. Usually, the vocabulary includes special tokens such as “BOS” (beginning of

sequence), “EOS” (end of sequence) and “PAD” (padding token, to ensure that the sequence is of length L).

4.8.7 Positional encoding

Due to their parameterization, Transformers are **equivariant** with respect to permutations: permuting the input sequence and applying the Transformer is equivalent to applying the Transformer and permuting the output sequence. Formally, for any permutation matrix \mathcal{P} of size $L \times L$ and input sequence $\mathbf{x}_1, \dots, \mathbf{x}_L$ seen as a matrix \mathbf{X} , we have

$$\text{Transformer}(\mathcal{P}\mathbf{X}) = \mathcal{P}\text{Transformer}(\mathbf{X}).$$

This means that Transformers treat an ordered sequence as a **multi-set** (a modification of the concept of a set that allows for multiple instances of its elements). To leverage order information in a Transformer, several approaches are possible (Dufter *et al.*, 2022): adding positional encoding (either absolute or relative), modifying the attention matrix and pre-processing the input sequence with an RNN. To add positional encoding, one typically adds a vector $\mathbf{p}_i \in \mathbb{R}^D$ representing the position $i \in [L]$ to the corresponding element $\mathbf{x}_i \in \mathbb{R}^D$,

$$\mathbf{x}_i \leftarrow \mathbf{x}_i + \mathbf{p}_i$$

An ideal positional encoding should work with any sequence length.

Learned positional encoding

Similarly to the token encoding, a simple way to define a positional encoding is to use an embedding layer. Each position $i \in [L]$ is transformed into a position vector $\mathbf{p}_i \in \mathbb{R}^D$ by

$$\mathbf{p}_i := \mathbf{W}^P \mathbf{e}_i$$

where $\mathbf{e}_i \in \mathbb{R}^L$ is the one-hot encoding of $i \in [L]$ and $\mathbf{W}^P \in \mathbb{R}^{D \times L}$ is a learnable embedding matrix. The column $\mathbf{W}_{:,i}^P \in \mathbb{R}^D$ can be interpreted as the continuous representation of position $i \in [L]$.

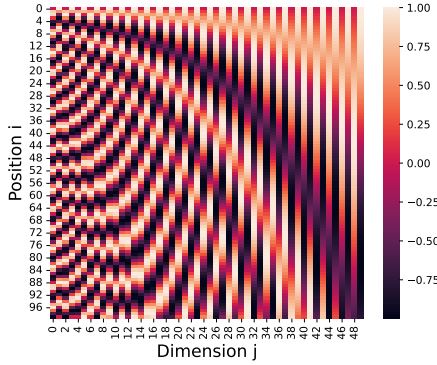


Figure 4.8: Heatmap of a positional encoding matrix $\mathbf{P} \in \mathbb{R}^{L \times D}$, with $L = 100$, $D = 50$, and $N = 30$. The values belong to the range $[-1, 1]$.

Sinusoidal positional encoding

Instead of learning the positional encoding, we can define it in a heuristic manner. For instance, Vaswani *et al.* (2017) proposed the **sinusoidal positional encoding** $\mathbf{p}_i := (p_{i,1}, \dots, p_{i,D})$, defined by

$$p_{i,j} := \begin{cases} \sin(\omega_j \cdot i) & \text{if } j \text{ is even} \\ \cos(\omega_{j-1} \cdot i) & \text{if } j \text{ is odd} \end{cases} \in [-1, 1],$$

where

$$\omega_j := \frac{1}{N^{\frac{j}{D}}} = N^{-\frac{j}{D}}$$

and where N is a constant, set to $N := 10000$ by the authors. We can gather the values in a position encoding matrix $\mathbf{P} \in \mathbb{R}^{L \times D}$, as depicted in Fig. 4.8.

On first sight, sinusoidal positional encoding may seem a bit mysterious. To gain some intuition, it is useful to compare it to a discrete

binary encoding of integers,

$$\begin{aligned}
 7 &\leftrightarrow 111 \\
 6 &\leftrightarrow 110 \\
 5 &\leftrightarrow 101 \\
 4 &\leftrightarrow 100 \\
 3 &\leftrightarrow 011 \\
 2 &\leftrightarrow 010 \\
 1 &\leftrightarrow 001 \\
 0 &\leftrightarrow 000.
 \end{aligned}$$

We see that the bits alternate more frequently between **0** and **1** as we go from right to left. The sinusoidal positional encoding achieves a similar behavior, but is continuous. Indeed, each coordinate $j \in [D]$ is associated with a sine wave (when j is even) or a cosine wave (when j is odd). The **wavelength** or **period** of the wave associated with j is $\frac{2\pi}{\omega_j}$ and its **frequency** is $\frac{\omega_j}{2\pi}$. Since ω_j is a decreasing function of j , we see that the waves oscillate more frequently when j is small. This behavior is illustrated in Fig. 4.9.

Stacking sine and cosine waves has been used for creating random Fourier features (Rahimi and Recht, 2007; Sutherland and Schneider, 2015). These use waves of random frequencies, while sinusoidal positional encoding uses waves of increasing frequency. These approaches therefore mainly differ in the way we choose wave frequencies.

Recovering relative positional information

An important property of the sinusoidal positional encoding is that $\mathbf{p}_{i+\delta}$ can be represented as a linear function of \mathbf{p}_i , for any offset δ (Vaswani *et al.*, 2017; Zhang *et al.*, 2021). The model should therefore be able to easily learn to attend by relative positions. Indeed, using the trigonometric identities for angle sums

$$\begin{aligned}
 \cos(\alpha + \beta) &= \cos(\alpha) \cos(\beta) - \sin(\alpha) \sin(\beta) \\
 \sin(\alpha + \beta) &= \sin(\alpha) \cos(\beta) + \cos(\alpha) \sin(\beta),
 \end{aligned}$$

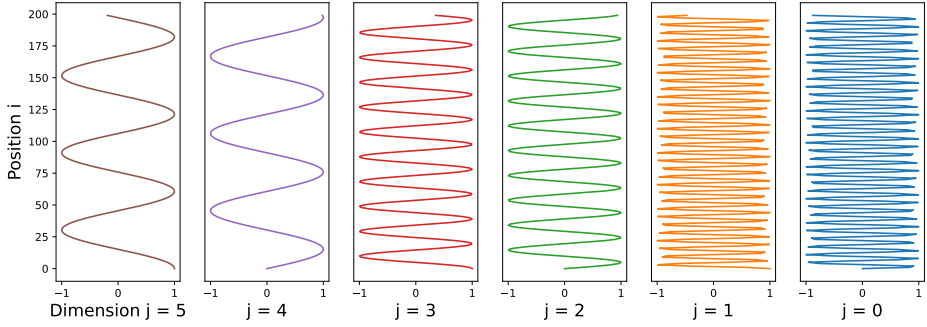


Figure 4.9: Using a sinusoidal positional encoding, each coordinate $j \in [D]$ is associated with a sine or cosine wave. In analogy with a binary encoding, in which lower bits alternate between 0 and 1 more frequently than higher bits, the wave associated with j oscillates more frequently when j is small.

we have for j even

$$\begin{aligned}
 \begin{pmatrix} p_{i+\delta,j+1} \\ p_{i+\delta,j} \end{pmatrix} &= \begin{pmatrix} \cos(\omega_j(i+\delta)) \\ \sin(\omega_j(i+\delta)) \end{pmatrix} \\
 &= \begin{pmatrix} \cos(\omega_j \cdot i) \cos(\omega_j \cdot \delta) - \sin(\omega_j \cdot i) \sin(\omega_j \cdot \delta) \\ \sin(\omega_j \cdot i) \cos(\omega_j \cdot \delta) + \cos(\omega_j \cdot i) \sin(\omega_j \cdot \delta) \end{pmatrix} \\
 &= \begin{pmatrix} \cos(\omega_j \cdot \delta) & -\sin(\omega_j \cdot \delta) \\ \sin(\omega_j \cdot \delta) & \cos(\omega_j \cdot \delta) \end{pmatrix} \begin{pmatrix} \cos(\omega_j \cdot i) \\ \sin(\omega_j \cdot i) \end{pmatrix} \\
 &= \begin{pmatrix} \cos(\omega_j \cdot \delta) & -\sin(\omega_j \cdot \delta) \\ \sin(\omega_j \cdot \delta) & \cos(\omega_j \cdot \delta) \end{pmatrix} \begin{pmatrix} p_{i,j+1} \\ p_{i,j} \end{pmatrix}
 \end{aligned}$$

Therefore, $(p_{i,j+1}, p_{i,j})$ can be linearly projected to $(p_{i+\delta,j+1}, p_{i+\delta,j})$, by applying a **rotation matrix**. This allows a Transformer to easily learn to attend by relative positions.

Rotary positional encoding (RoPE)

Another popular approach for taking into account absolute positional information in self-attention is RoPE (Su *et al.*, 2024), which stands for rotary positional encoding. For simplicity, we briefly explain the idea with a single attention head and $D_v = D_k = D$. Suppose we have a sequence $\mathbf{x}_1, \dots, \mathbf{x}_L \in \mathbb{R}^D$ already encoded with token encoding. We

then define the query and key vectors as

$$\begin{aligned}\mathbf{q}_m &:= \mathbf{R}_D(\theta, m) \mathbf{W}^Q \mathbf{x}_m \quad m \in [L] \\ \mathbf{k}_n &:= \mathbf{R}_D(\theta, n) \mathbf{W}^K \mathbf{x}_n \quad n \in [L],\end{aligned}$$

where $\mathbf{R}_D(\theta, m) \in \mathbb{R}^{D \times D}$ is a rotation matrix. The main intuition is that we are rotating the 2-dimensional vector by an angle which is a multiple of the position m . In contrast to the learned and sinusoidal positional encodings, the weight matrices \mathbf{W}^Q and \mathbf{W}^K are applied **before** applying RoPE, not after.

When $D = 2$, $\mathbf{R}_D(\theta, m) = \mathbf{R}_2(\theta, m)$ is defined as the rotation matrix

$$\mathbf{R}_2(\theta, m) := \begin{pmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{pmatrix}$$

for some angle $\theta \in \mathbb{R}$ and position $m \in [L]$. Using Euler's formula

$$e^{im\theta} = \cos(m\theta) + i \sin(m\theta),$$

we note that $\mathbf{v}' := \mathbf{R}_2(\theta, m) \mathbf{v}$ for $\mathbf{v} = (v_1, v_2) \in \mathbb{R}^2$ and $\mathbf{v}' = (v'_1, v'_2) \in \mathbb{R}^2$ is equivalent to $z' := z e^{im\theta}$ for $z := v_1 + iv_2 \in \mathbb{C}$ and $z' = v'_1 + iv'_2 \in \mathbb{C}$. See Su *et al.* (2024) for a detailed derivation of RoPE's formula.

When $D > 2$, assuming D is even, $\mathbf{R}_D(\theta, m)$ is defined as

$$\mathbf{R}_D(\theta, m) := \begin{pmatrix} \mathbf{R}_2(\theta_1, m) & & & \\ & \mathbf{R}_2(\theta_2, m) & & \\ & & \ddots & \\ & & & \mathbf{R}_2(\theta_{D/2}, m) \end{pmatrix},$$

where $\theta_j := 10000^{-2(j-1)/D}$. In practice, we never materialize $\mathbf{R}_D(\theta, m)$ as a matrix, since it would be very sparse, but rather view it as a linear map applied to an arbitrary vector $\mathbf{v} \in \mathbb{R}^D$,

$$\mathbf{R}_D(\theta, m) \mathbf{v} = \begin{pmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \\ \vdots \\ v_{D-1} \\ v_D \end{pmatrix} \otimes \begin{pmatrix} \cos(m\theta_1) \\ \cos(m\theta_1) \\ \cos(m\theta_2) \\ \cos(m\theta_2) \\ \vdots \\ \cos(m\theta_{D/2}) \\ \cos(m\theta_{D/2}) \end{pmatrix} + \begin{pmatrix} -v_2 \\ v_1 \\ -v_4 \\ v_3 \\ \vdots \\ -v_D \\ v_{D-1} \end{pmatrix} \otimes \begin{pmatrix} \sin(m\theta_1) \\ \sin(m\theta_1) \\ \sin(m\theta_2) \\ \sin(m\theta_2) \\ \vdots \\ \sin(m\theta_{D/2}) \\ \sin(m\theta_{D/2}) \end{pmatrix}.$$

Once we computed \mathbf{Q} and \mathbf{K} , we can use $\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}; \mathbf{M})$ as usual. For multi-head attention, we simply apply the same approach for each attention head.

As analyzed by Su *et al.* (2024), RoPE satisfies several valuable properties. In particular, it is flexible w.r.t. the sequence length L and despite being an absolute encoding, it manages to also capture relative distance between tokens. It therefore has the merits of both absolute and relative positional encodings.

4.8.8 Decoder-only architectures

Defining a language model

Transformers were originally developed to create encoder-decoder architectures for machine translation (Vaswani *et al.*, 2017). However, in the context of language modelling, Transformers are now routinely used for decoder-only **autoregressive** architectures. Suppose we are given a sequence of discrete symbols $x_1, \dots, x_L \in [M]$ (if a sequence has less than L elements, we can use padding symbols). The goal of (unconditional) language models is to create a function producing the joint probability,

$$(x_1, \dots, x_L) \mapsto \mathbb{P}(X_1 = x_1, \dots, X_L = x_L).$$

Using the chain rule of probability (Section 10.1), without loss of generality, we can write

$$\mathbb{P}(X_1 = x_1, \dots, X_L = x_L) = \prod_{k=1}^L \mathbb{P}(X_k = x_k \mid X_1 = x_1, \dots, X_{k-1} = x_{k-1}).$$

If the maximum length L is very large, it is more practical to only consider a context window of size N ,

$$\mathbb{P}(X_1 = x_1, \dots, X_L = x_L) := \prod_{k=1}^L \mathbb{P}(X_k = x_k \mid X_{k-N} = x_{k-N}, \dots, X_{k-1} = x_{k-1}).$$

This amounts to defining a **higher-order Markov chain** (Section 10.4.3). Using this approach, creating a language model then boils down to defining a function

$$(x_{k-N}, \dots, x_{k-1}, x_k) \mapsto \mathbb{P}(X_k = x_k \mid X_{k-N} = x_{k-N}, \dots, X_{k-1} = x_{k-1}).$$

Let us assume that the sequence of discrete symbols $x_1, \dots, x_L \in [M]$ has been mapped to a sequence of continuous vectors $\mathbf{x}_1, \dots, \mathbf{x}_L \in \mathbb{R}^D$ using token encoding (Section 4.8.6) and positional encoding (Section 4.8.7). Using a causal Transformer block, we can obtain a representation of the current context,

$$(\mathbf{x}_{k-N}, \dots, \mathbf{x}_{k-1}) \leftarrow \text{Transformer}(\mathbf{x}_{k-N}, \dots, \mathbf{x}_{k-1}; \mathbf{M}_N).$$

To ensure that the Transformer only relies on past tokens to predict the current token, the mask matrix \mathbf{M}_N is set to a **triangular matrix** of size $N \times N$, with the elements of the lower part set to 1, and the elements of the upper part set to $-\infty$. To reduce this to logits in \mathbb{R}^M , we usually use the last element of the context $\mathbf{x}_{k-1} \in \mathbb{R}^D$ and apply a linear model, to obtain

$$\boldsymbol{\theta}_k := \mathbf{W}^D \mathbf{x}_{k-1} \in \mathbb{R}^M,$$

where $\mathbf{W}^D \in \mathbb{R}^{M \times D}$ is a learned “disembedding” matrix. To obtain a valid probability distribution, we apply a soft-argmax on the logits $\boldsymbol{\theta}_k$,

$$\boldsymbol{\pi}_k := \text{softargmax}(\boldsymbol{\theta}_k) \in \Delta^M.$$

Finally, we can now define

$$\mathbb{P}(X_k = x_k \mid X_{k-N} = x_{k-N}, \dots, X_{k-1} = x_{k-1}) := [\boldsymbol{\pi}_k]_{x_k}.$$

Remark 4.5 (Weight tying). Instead of learning a separate disembedding matrix $\mathbf{W}^D \in \mathbb{R}^{M \times D}$, a frequently used technique is to set $\mathbf{W}^D := (\mathbf{W}^E)^\top$, where $\mathbf{W}^E \in \mathbb{R}^{D \times M}$ is the embedding matrix from Section 4.8.6. This weight tying reduces the number of parameters to learn and is shown to work well in practice (Press and Wolf, 2016).

Training

Let us gather the Transformer parameters (multi-head attention, MLP, layer norm) as $\mathbf{w} \in \mathbb{R}^P$. The decoder-only Transformer with a context window of size N then defines a parametric probability distribution

$$p_{\mathbf{w}}(x_1, \dots, x_L) := \prod_{k=1}^L p_{\mathbf{w}}(x_k \mid x_{k-N}, \dots, x_{k-1}).$$

Given a corpus of sequences \mathcal{D} , we usually seek the parameters $\mathbf{w} \in \mathbb{R}^P$ by maximizing the log-likelihood,

$$\mathbb{E}_{x \sim \mathcal{D}} [\log p_{\mathbf{w}}(x_1, \dots, x_L)] = \mathbb{E}_{x \sim \mathcal{D}} \left[\sum_{k=1}^L \log p_{\mathbf{w}}(x_k \mid x_{k-N}, \dots, x_{k-1}) \right].$$

This amounts to using a logistic (cross-entropy) loss in a token-wise fashion. Importantly, the context x_{k-N}, \dots, x_{k-1} used to predict the next token x_k always comes from the data, not from the tokens generated by the model. This is called **teacher forcing** and makes Transformer training highly parallelizable. The objective is usually solved approximately using stochastic gradient algorithms. To maximize GPU utilization, sequences of variable lengths are usually packed together in order to form batches.

Sampling

A decoder-only Transformer with a context window of size N can be seen as forming a higher-order Markov chain, which is a special case of Bayesian network. To generate a sequence from the model, we can therefore use **ancestral sampling** (Section 10.5.3),

$$\begin{aligned} X_0 &:= x_0 \\ X_1 &\sim p_{\mathbf{w}}(\cdot \mid X_0) \\ X_2 &\sim p_{\mathbf{w}}(\cdot \mid X_0, X_1) \\ &\vdots \\ X_k &\sim p_{\mathbf{w}}(\cdot \mid X_{k-N}, \dots, X_{k-1}), \end{aligned}$$

where x_0 denotes the beginning-of-sequence (BOS) token. The sampling stops when the end-of-sequence (EOS) token is generated or when a maximum length is reached. The generated sequences are i.i.d.

Because sampling happens one token at a time, it is highly sequential.

The Transformer block is called sequentially as

$$\begin{aligned} &\text{Transformer}(X_0; \mathbf{M}_1) \\ &\text{Transformer}(X_0, X_1; \mathbf{M}_2) \\ &\dots \\ &\text{Transformer}(X_{k-1}, \dots X_{k-N}; \mathbf{M}_N). \end{aligned}$$

To avoid repeating the same computations again and again, assuming that a causal mask is used, the past key and value matrices used in multi-head attention (Section 4.8.3) are usually stored in the so-called **KV cache**.

4.8.9 Encoder-only architectures

Encoder-only Transformers can be used for learning to represent sequences. The most prominent example is BERT (Devlin *et al.*, 2019), which stands for bidirectional encoder representations from transformers. BERT uses Algorithm 4.4 without causal mask (hence “bidirectional” in the acronym). Discrete sequence elements are transformed to vectors using token encoding, positional encoding and potentially segment encoding. For reasons that will become clear below, the first token of every sequence is always a special classification token [CLS]. Training is broken down into two phases: pretraining and finetuning.

Pretraining

Since pre-training is performed on an unlabeled corpus, it is necessary to synthetically generate prediction tasks.

In masked prediction, a percentage (typically 15%) of the input tokens are randomly masked. The model is then trained to predict the original vocabulary ID of the masked tokens, given the context provided by the unmasked tokens. This forces the model to learn a rich, bidirectional representation of the language. The masking strategy involves:

- 80% of the time, replacing the chosen token with [MASK],
- 10% of the time, replacing the chosen token with a random token from the vocabulary,

- 10% of the time, keeping the chosen token unchanged.

In next sentence prediction, the model is trained to understand the relationship between two sentences. For each training example, BERT is given two sentences, A and B . 50% of the time, B is the actual next sentence that follows A in the original document (labeled *IsNext*). The other 50% of the time, B is a random sentence from the corpus (labeled *NotNext*). A special [CLS] token is prepended to the input sequence, and its final hidden state is used to predict whether sentence B follows sentence A . The [SEP] token separates the two sentences.

Finetuning

After pre-training on a massive corpus, the pretrained BERT model can be finetuned for various downstream NLP tasks (e.g., text classification, question answering, named entity recognition) by adding a small, task-specific output layer on top of the pre-trained Transformer encoder. The entire model, including the pretrained weights, is then finetuned on the labeled data for the specific task. For classification tasks, the final hidden state corresponding to the class token [CLS] is used. For tagging tasks (token-level classification) such as named entity recognition (NER) or part-of-speech (POS) tagging, the final hidden state of each token is used. For question answering tasks, which usually involve finding a span (start and end token) within a given text, this is treated as two token-level classification problems: one for the start token and one for the end token.

4.8.10 Encoder-decoder architectures

The encoder-decoder architecture was the original architecture proposed in the seminal Transformer paper (Vaswani *et al.*, 2017). It can be used for sequence-to-sequence tasks, such as machine translation, summarization and question answering. We denote the input sequence by \mathbf{X} and the output sequence by \mathbf{Y} .

Encoder

The role of the encoder is to produce a rich representation \mathbf{H}_{enc} of the input sequence \mathbf{X} , which serves as a context for the subsequent decoding phase. As in encoder-only architectures, this is achieved by using Algorithm 4.4 without causal mask. The multi-head self-attention layer allows the encoder to weigh the importance of different tokens in the input sequence relative to each other. For each token in the input, it computes a representation that incorporates information from all other tokens in the same sequence. The encoder block uses its own parameters for each sub-component (multi-head attention, layer norm, MLP). In particular, multi-head attention uses parameters $\{\mathbf{W}_i^Q\}_{i=1}^H$, $\{\mathbf{W}_i^K\}_{i=1}^H$, $\{\mathbf{W}_i^V\}_{i=1}^H$ and \mathbf{W}^O .

Decoder

The role of the decoder is to generate the output sequence \mathbf{Y} one token at a time, based on the encoded representation \mathbf{H}_{enc} from the encoder and the previously generated tokens $\mathbf{Y}_{1:t-1}$. The decoder differs in the way attention is applied. First, we apply causal multi-head self-attention to obtain a representation \mathbf{H}_{dec} corresponding to the previously generated tokens $\mathbf{Y}_{1:t-1}$. Second, we apply multi-head **cross-attention**. The key difference with multi-head self-attention is that the key, query and value matrices are defined as

$$\begin{aligned} \mathbf{Q}_i^{\text{dec}} &:= \mathbf{H}_{\text{dec}} \mathbf{W}_i'^Q \\ \mathbf{K}_i^{\text{enc}} &:= \mathbf{H}_{\text{enc}} \mathbf{W}_i'^K \\ \mathbf{V}_i^{\text{enc}} &:= \mathbf{H}_{\text{enc}} \mathbf{W}_i'^V, \end{aligned}$$

where $\{\mathbf{W}_i'^Q\}_{i=1}^H$, $\{\mathbf{W}_i'^K\}_{i=1}^H$ and $\{\mathbf{W}_i'^V\}_{i=1}^H$ are the weight matrices used for the decoder. Subsequently, MLP and layer norm layers are applied, similarly as before.

Differences with decoder-only architectures

An encoder-decoder architecture provides a dedicated component for understanding the input and another for generating the output. It allows

for a bidirectional understanding of the input. It is therefore best suited for sequence-to-sequence (seq2seq) tasks where the input and output sequences may have different domains or modalities.

Decoder-only models, on the other hand, are streamlined for generating text by continuously predicting the next token, leveraging a single, powerful autoregressive mechanism. That being said, due to their simplicity, decoder-only architectures are now increasingly being used even for seq2seq tasks, by prepending the input (prompt) to the context. This means that the input is subject to the causal mask, but the performance remains remarkably strong.

4.9 Summary

- Programs can be mathematically represented as a directed acyclic graph.
- Neural networks are parameterized programs.
- Feedforward networks are parameterized computation chains.
- Multilayer perceptrons (MLPs), residual neural networks (ResNets) and convolutional neural network (CNNs) are all particular parametrizations of feedforward networks.
- Transformer blocks are designed to process sequences of variable-length, but they are equivariant to permutations and therefore require positional encoding, in order to leverage positional information.

5

Control flows

Control flows, such as conditionals or loops, are an essential part of computer programming, as they allow us to express complex programs. It is therefore natural to ask whether these constructs can be included in a differentiable program. This is what we study in this chapter.

5.1 Comparison operators

Control flows rely on **comparison operators**, a.k.a. **relational operators**. Formally, we can define a comparison operator $\pi = \text{op}(u_1, u_2)$ as a function from $u_1 \in \mathbb{R}$ and $u_2 \in \mathbb{R}$ to $\pi \in \{0, 1\}$. The binary (Boolean) output π can then be used within a conditional statement (see Section 5.6, Section 5.7) to decide whether to execute one branch or another. We define the following operators, illustrated in Fig. 5.1:

- **greater than:**

$$\begin{aligned}\text{gt}(u_1, u_2) &:= \begin{cases} 1 & \text{if } u_1 \geq u_2 \\ 0 & \text{otherwise} \end{cases} \\ &= \text{step}(u_1 - u_2)\end{aligned}$$

- **less than:**

$$\begin{aligned} \text{lt}(u_1, u_2) &:= \begin{cases} 1 & \text{if } u_1 \leq u_2 \\ 0 & \text{otherwise} \end{cases} \\ &= 1 - \text{gt}(u_1, u_2) \\ &= \text{step}(u_2 - u_1) \end{aligned}$$

- **equal:**

$$\begin{aligned} \text{eq}(u_1, u_2) &:= \begin{cases} 1 & \text{if } |u_1 - u_2| = 0 \\ 0 & \text{otherwise} \end{cases} \\ &= \text{gt}(u_2, u_1) \cdot \text{gt}(u_1, u_2) \\ &= \text{step}(u_2 - u_1) \cdot \text{step}(u_1 - u_2) \end{aligned}$$

- **not equal:**

$$\begin{aligned} \text{neq}(u_1, u_2) &:= \begin{cases} 1 & \text{if } |u_1 - u_2| > 0 \\ 0 & \text{otherwise} \end{cases} \\ &= 1 - \text{eq}(u_1, u_2) \\ &= 1 - \text{step}(u_2 - u_1) \cdot \text{step}(u_1 - u_2), \end{aligned}$$

where $\text{step}: \mathbb{R} \rightarrow \{0, 1\}$ is the **Heaviside step function**

$$\text{step}(u) := \begin{cases} 1 & \text{if } u \geq 0 \\ 0 & \text{otherwise} \end{cases}.$$

The Heaviside step function is piecewise constant. At $u = 0$, the function is discontinuous. At $u \neq 0$, it is continuous and has **null derivative**. Since the comparison operators we presented are all expressed in terms of the step function, they are all continuous and differentiable almost everywhere, with null derivative. Therefore, while their derivatives are well-defined almost everywhere, they are **uninformative** and prevent gradient backpropagation.

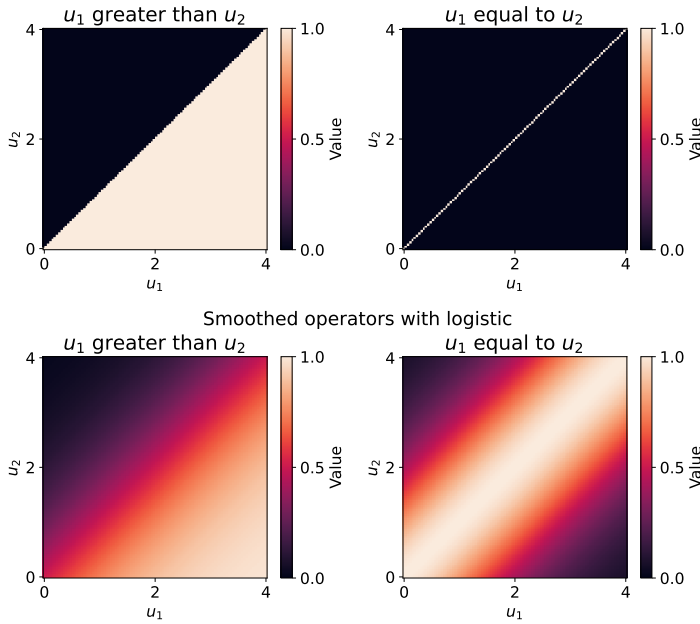


Figure 5.1: The greater than and equal to operators are discontinuous functions, leading to black or white pictures. They can be smoothed with appropriate approximations of the Heaviside step function.

5.2 Soft inequality operators

5.2.1 Heuristic definition

To obtain a **continuous relaxation** of inequality operators, we can heuristically replace the step function in the expression of “greater than” and “less than” by a sigmoid function sigmoid_σ , where $\sigma > 0$ is a scaling parameter. Such a sigmoid function should satisfy the following properties:

- $\text{sigmoid}_\sigma(-u) = 1 - \text{sigmoid}_\sigma(u)$,
- $\lim_{u \rightarrow \infty} \text{sigmoid}_\sigma(u) = 1$,
- $\lim_{u \rightarrow -\infty} \text{sigmoid}_\sigma(u) = 0$,
- $\text{sigmoid}_\sigma(0) = \frac{1}{2}$.

Two examples of sigmoids satisfying the aforementioned properties are the logistic function (the CDF of the standard logistic distribution)

$$\text{sigmoid}_\sigma(u) := \text{logistic}_\sigma(u) := \frac{1}{1 + e^{-u/\sigma}} \in (0, 1)$$

and the standard Gaussian's CDF, defined in Eq. (3.1),

$$\text{sigmoid}_\sigma(u) := \Phi(u/\sigma).$$

We may then define the soft “greater than”

$$\begin{aligned} \text{gt}(\mu_1, \mu_2) &= \text{step}(\mu_1 - \mu_2) \\ &\approx \text{sigmoid}_\sigma(\mu_1 - \mu_2) \\ &=: \text{gt}_\sigma(\mu_1, \mu_2) \end{aligned}$$

and the soft “less than”

$$\begin{aligned} \text{lt}(\mu_1, \mu_2) &= \text{step}(\mu_2 - \mu_1) \\ &\approx \text{sigmoid}_\sigma(\mu_2 - \mu_1) \\ &=: \text{lt}_\sigma(\mu_1, \mu_2) \\ &= 1 - \text{sigmoid}_\sigma(\mu_1 - \mu_2) \\ &= 1 - \text{gt}_\sigma(\mu_1 - \mu_2). \end{aligned}$$

In the limit, we have that $\text{sigmoid}_\sigma(\mu_1 - \mu_2) \rightarrow 1$ when $\mu_1 - \mu_2 \rightarrow \infty$. In the limit, sigmoid_σ therefore outputs a value of 1 if μ_1 and μ_2 are infinitely apart. Besides the logistic function and the standard Gaussian's CDF, other sigmoid functions are possible, as discussed in Section 13.6. In particular, with sparse sigmoids, there exists a finite value τ such that $\mu_1 - \mu_2 \geq \tau \implies \text{sigmoid}_\sigma(\mu_1 - \mu_2) = 1$.

5.2.2 Stochastic process perspective

When the sigmoid used to replace the step function is the logistic function or the standard Gaussian's CDF, we can revisit the previous heuristic definition of $\text{gt}_\sigma(\mu_1, \mu_2)$ and $\text{lt}_\sigma(\mu_1, \mu_2)$ from a more formal perspective. Indeed, to real values $\mu_1 \in \mathbb{R}$ and $\mu_2 \in \mathbb{R}$, we can associate random variables

$$U_1 \sim p_{\mu_1, \sigma_1} \quad \text{and} \quad U_2 \sim p_{\mu_2, \sigma_2},$$

thereby forming a **stochastic process** (we assume that σ_1 and σ_2 are fixed). Alternatively, we can also write

$$(U_1, U_2) \sim p_{\mu_1, \sigma_1} \otimes p_{\mu_2, \sigma_2},$$

where for two distributions p_1 and p_2 , we denote $p_1 \otimes p_2$ their outer product $(p_1 \otimes p_2)(u_1, u_2) := p_1(u_1)p_2(u_2)$. We can then define

$$\begin{aligned} \text{gt}_\sigma(\mu_1, \mu_2) &= \mathbb{E} [\text{gt}(U_1, U_2)] \\ &= \mathbb{E} [\text{step}(U_1 - U_2)] \\ &= \mathbb{P}(U_1 - U_2 \geq 0) \\ &= \mathbb{P}(U_2 - U_1 \leq 0) \\ &= F_{U_2 - U_1}(0), \end{aligned}$$

where F_X is the **cumulative distribution function** (CDF) of the random variable X , and σ is a function of σ_1 and σ_2 . Similarly, we obtain

$$\begin{aligned} \text{lt}_\sigma(\mu_1, \mu_2) &= \mathbb{E} [\text{lt}(U_1, U_2)] \\ &= \mathbb{E} [\text{step}(U_2 - U_1)] \\ &= \mathbb{P}(U_1 - U_2 \leq 0) \\ &= F_{U_1 - U_2}(0). \end{aligned}$$

We see that the soft inequality operators are based on the CDF of the **difference** between U_1 and U_2 .

For location-scale family distributions (Section 12.4.1), from a perturbation perspective, we can also define noise variables $Z_1 \sim p_{0,1}$ and $Z_2 \sim p_{0,1}$ such that $U_1 = \mu_1 + \sigma_1 Z_1$ and $U_2 = \mu_2 + \sigma_2 Z_2$ (Section 12.4.1). We then have

$$\begin{aligned} \text{gt}_\sigma(\mu_1, \mu_2) &= \mathbb{E} [\text{gt}(\mu_1 + \sigma_1 Z_1, \mu_2 + \sigma_2 Z_2)] \\ \text{lt}_\sigma(\mu_1, \mu_2) &= \mathbb{E} [\text{lt}(\mu_1 + \sigma_1 Z_1, \mu_2 + \sigma_2 Z_2)]. \end{aligned}$$

Gaussian case

When $U_1 \sim \text{Normal}(\mu_1, \sigma_1^2)$ and $U_2 \sim \text{Normal}(\mu_2, \sigma_2^2)$, we have

$$U_1 - U_2 \sim \text{Normal}(\mu_1 - \mu_2, \sigma_1^2 + \sigma_2^2). \quad (5.1)$$

Denoting Φ the standard Gaussian's CDF, we then obtain

$$\begin{aligned}
 \text{gt}_\sigma(\mu_1, \mu_2) &= F_{U_2 - U_1}(0) \\
 &= \Phi\left(\frac{0 - (\mu_2 - \mu_1)}{\sigma}\right) \\
 &= \Phi\left(\frac{\mu_1 - \mu_2}{\sigma}\right) \\
 \text{lt}_\sigma(\mu_1, \mu_2) &= F_{U_1 - U_2}(0) \\
 &= \Phi\left(\frac{0 - (\mu_1 - \mu_2)}{\sigma}\right) \\
 &= \Phi\left(\frac{\mu_2 - \mu_1}{\sigma}\right),
 \end{aligned}$$

where $\sigma := \sqrt{\sigma_1^2 + \sigma_2^2}$. The corresponding distribution for Z_1 and Z_2 is **Gaussian noise**.

Logistic case

When $U_1 \sim \text{Gumbel}(\mu_1, \sigma)$ and $U_2 \sim \text{Gumbel}(\mu_2, \sigma)$, we have

$$U_1 - U_2 \sim \text{Logistic}(\mu_1 - \mu_2, \sigma). \quad (5.2)$$

We then obtain (see also Proposition 14.3)

$$\begin{aligned}
 \text{gt}_\sigma(\mu_1, \mu_2) &= \text{logistic}\left(\frac{\mu_1 - \mu_2}{\sigma}\right) \\
 \text{lt}_\sigma(\mu_1, \mu_2) &= \text{logistic}\left(\frac{\mu_2 - \mu_1}{\sigma}\right).
 \end{aligned}$$

The corresponding distribution for Z_1 and Z_2 is **Gumbel noise**.

Recovering hard inequality operators

We easily recover the “hard” inequality operator by

$$\text{gt}(\mu_1, \mu_2) = \mathbb{E}[\text{gt}(U_1, U_2)],$$

where $U_i \sim \delta_{\mu_i}$ and where δ_{μ_i} is the **delta distribution** that assigns a probability of 1 to μ_i .

5.3 Soft equality operators

5.3.1 Heuristic definition

The equality operator $\text{eq}(\mu_1, \mu_2)$ can be seen as an extreme kind of **similarity function** between numbers, that can only output the values 0 or 1. To define soft equality operators, a natural idea is therefore to replace the equality operator by a more general similarity function. A similarity function should achieve its maximum at $\mu_1 = \mu_2$ and it should decrease as μ_1 and μ_2 move apart. A common family of similarity functions are **kernels**. Briefly, a kernel $k(\mu_1, \mu_2)$ can be seen as the inner product

$$k(\mu_1, \mu_2) := \langle \phi(\mu_1), \phi(\mu_2) \rangle$$

between the embeddings $\phi(\mu_1)$ and $\phi(\mu_2)$ of μ_1 and μ_2 in some (potentially infinite-dimensional) space \mathcal{H} , a reproducing kernel Hilbert space to be precise; see Schölkopf and Smola (2002) and Shawe-Taylor and Cristianini (2004) for an in-depth review of kernels. To obtain a similarity measure between 0 and 1 approximating the equality operator, we can normalize to obtain

$$\begin{aligned} \text{eq}(\mu_1, \mu_2) &\approx \frac{k(\mu_1, \mu_2)}{\sqrt{k(\mu_1, \mu_1)k(\mu_2, \mu_2)}} \\ &= \frac{\langle \phi(\mu_1), \phi(\mu_2) \rangle}{\|\phi(\mu_1)\| \|\phi(\mu_2)\|}, \end{aligned}$$

where $\|\phi(\mu)\| := \sqrt{\langle \phi(\mu), \phi(\mu) \rangle} = \sqrt{\kappa(\mu, \mu)}$. This is the **cosine similarity** between $\phi(\mu_1)$ and $\phi(\mu_2)$.

A particular kind of kernel are **isotropic** kernels of the form

$$k(\mu_1, \mu_2) := \kappa(\mu_1 - \mu_2),$$

that depend only on the difference between inputs. When the kernel has a scale parameter $\sigma > 0$, we use the notation κ_σ . We can then define a soft equality operator as

$$\text{eq}(\mu_1, \mu_2) \approx \text{eq}_\sigma(\mu_1, \mu_2) := \frac{\kappa_\sigma(\mu_1 - \mu_2)}{\kappa_\sigma(0)}.$$

Several isotropic kernels can be chosen such as the **Gaussian kernel**

$$\kappa_\sigma(t) := \exp\left(-\frac{t^2}{2\sigma^2}\right)$$

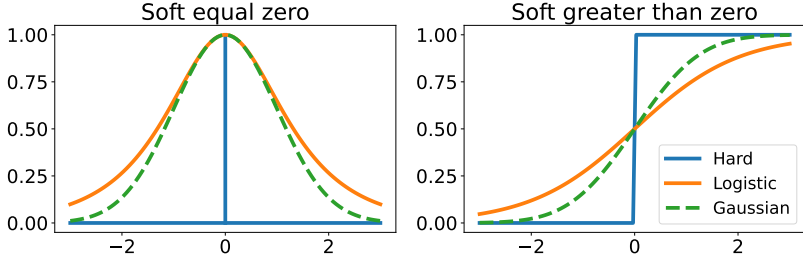


Figure 5.2: Soft equality and soft greater than operators can be defined as normalized kernels (PDF) and as CDF functions, respectively.

or the **logistic kernel**

$$\kappa_{\sigma}(t) := \text{sech}^2\left(\frac{t}{2\sigma}\right),$$

where we defined the hyperbolic secant

$$\text{sech}(u) := 2/(\exp(u) + \exp(-u)).$$

As their names suggest, these kernels arise naturally from a probabilistic perspective, that we present below.

The soft equality operators obtained with these kernels are illustrated in Fig. 5.2. Intuitively, we replaced a bar located at $\mu_1 = \mu_2$ with a bump function. The soft equality operator obtained with the logistic kernel coincides with the expression Petersen *et al.* (2021) arrive at (see their Eq. 9), in a different manner.

5.3.2 Stochastic process perspective

We again adopt the stochastic process perspective, in which we associate random variables

$$U_1 \sim p_{\mu_1, \sigma_1} \quad \text{and} \quad U_2 \sim p_{\mu_2, \sigma_2}$$

to real values $\mu_1 \in \mathbb{R}$ and $\mu_2 \in \mathbb{R}$. However, to handle the equality operator, we cannot simply use the expectation of $\text{eq}(U_1, U_2)$ since

$$\mathbb{E}[\text{eq}(U_1, U_2)] = \mathbb{P}(U_1 = U_2) = 0,$$

U_1 and U_2 being independent continuous variables. While we cannot use the probability of $U_1 = U_2$, or equivalently of $U_1 - U_2 = 0$, we can

consider using the probability density function (PDF) $f_{U_1-U_2}$ of $U_1 - U_2$ evaluated at 0. To ensure that the maximum is achieved at 0 with value 1, we can normalize the PDF to define

$$\text{eq}_\sigma(\mu_1, \mu_2) = \frac{f_{U_1-U_2}(0)}{f_0(0)}.$$

It is well-known that the PDF of the sum of two random variables is the **convolution** of their respective PDFs. We therefore have

$$\begin{aligned} f_{U_1-U_2}(t) &= (f_{U_1} * f_{-U_2})(t) \\ &= \int_{-\infty}^{\infty} f_{U_1}(\tau) f_{-U_2}(t - \tau) d\tau. \end{aligned}$$

In particular, with $t = 0$, if f_X is the PDF of a location-scale family distributed random variable, we obtain

$$\begin{aligned} f_{U_1-U_2}(0) &= (f_{U_1} * f_{-U_2})(0) \\ &= \int_{-\infty}^{\infty} f_{U_1}(\tau) f_{-U_2}(-\tau) d\tau \\ &= \int_{-\infty}^{\infty} f_{U_1}(\tau) f_{U_2}(\tau) d\tau \\ &:= \langle f_{U_1}, f_{U_2} \rangle \\ &:= \kappa(\mu_1, \mu_2). \end{aligned}$$

We indeed recover an inner product and therefore a kernel.

CDF and PDF of absolute difference

While $\mathbb{P}(U_1 = U_2) = 0$, we can also consider using $\mathbb{P}(|U_1 - U_2| \leq \varepsilon) = F_{|U_1-U_2|}(\varepsilon)$ as an alternative notion of soft equality. For any random variable X , we have

$$\begin{aligned} F_{|X|}(x) &= \mathbb{P}(|X| \leq x) \\ &= \mathbb{P}(-x \leq X \leq x) \\ &= \mathbb{P}(X \leq x) - \mathbb{P}(X \leq -x) \\ &= F_X(x) - F_X(-x). \end{aligned}$$

Therefore,

$$\mathbb{P}(|U_1 - U_2| \leq \varepsilon) = F_{U_1-U_2}(\varepsilon) - F_{U_1-U_2}(-\varepsilon).$$

We can also derive the PDF of $|X|$ as

$$\begin{aligned} f_{|X|}(x) &= F'_X(x) - F'_X(-x) \\ &= f_X(x) + f_X(-x) \end{aligned}$$

and in particular

$$f_{|X|}(0) = 2f_X(0).$$

Therefore

$$f_{|U_1 - U_2|}(0) = 2f_{U_1 - U_2}(0),$$

further justifying using the PDF of $U_1 - U_2$ evaluated at 0. When X follows a normal distribution, $|X|$ follows the so-called folded normal distribution.

Gaussian case

When $U_1 \sim \text{Normal}(\mu_1, \sigma_1^2)$ and $U_2 \sim \text{Normal}(\mu_2, \sigma_2^2)$, we obtain from Eq. (5.1)

$$f_{U_1 - U_2}(t) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(t - (\mu_1 - \mu_2))^2}{2(\sigma_1^2 + \sigma_2^2)}\right)$$

so that

$$\text{eq}_\sigma(\mu_1, \mu_2) = \exp\left(\frac{(\mu_1 - \mu_2)^2}{2(\sigma_1^2 + \sigma_2^2)}\right) \in [0, 1].$$

We indeed recover $\kappa_\sigma(\mu_1 - \mu_2)/\kappa_\sigma(0)$, where κ_σ is the Gaussian kernel with $\sigma = \sqrt{\sigma_1^2 + \sigma_2^2}$. For the CDF of the absolute difference, we obtain

$$\mathbb{P}(|U_1 - U_2| \leq \varepsilon) = \Phi\left(\frac{\varepsilon - (\mu_1 - \mu_2)}{\sigma}\right) - \Phi\left(\frac{-\varepsilon - (\mu_1 - \mu_2)}{\sigma}\right).$$

Logistic case

When $U_1 \sim \text{Gumbel}(\mu_1, \sigma)$ and $U_2 \sim \text{Gumbel}(\mu_2, \sigma)$, recalling that

$$\text{sech}(u) := 2/(\exp(u) + \exp(-u)),$$

we obtain from Eq. (5.2)

$$f_{U_1 - U_2}(t) = \frac{1}{4\sigma} \text{sech}^2\left(\frac{t - (\mu_1 - \mu_2)}{2\sigma}\right)$$

so that

$$\text{eq}_\sigma(\mu_1, \mu_2) = \text{sech}^2\left(\frac{\mu_1 - \mu_2}{2\sigma}\right) \in [0, 1].$$

We indeed recover $\kappa_\sigma(\mu_1 - \mu_2)/\kappa_\sigma(0)$, where κ_σ is the logistic kernel with $\sigma = \sigma_1 = \sigma_2$.

5.3.3 Gaussian process perspective

The previous approach relied on mapping μ_1 and μ_2 to two **independent** random variables $U_1 \sim p_{\mu_1, \sigma_1}$ and $U_2 \sim p_{\mu_2, \sigma_2}$ (we assume that σ_1 and σ_2 are fixed). Instead, we can consider mapping μ_1 and μ_2 to two **dependent** random variables U_1 and U_2 , whose covariance depends on the similarity between μ_1 and μ_2 . We can do so by using a **Gaussian process** (Hida and Hitsuda, 1976).

A Gaussian process on \mathbb{R} is a stochastic process $\{U_\mu : \mu \in \mathbb{R}\}$ indexed by $\mu \in \mathbb{R}$ such that any subset of K random variables $(U_{\mu_1}, \dots, U_{\mu_K})$ associated with $(\mu_1, \dots, \mu_K) \in \mathbb{R}$ is a multivariate Gaussian random variable. The Gaussian process is characterized by the mean function $\mu \mapsto \mathbb{E}[U_\mu]$, and its covariance function $(\mu_i, \mu_j) \mapsto \text{Cov}(U_{\mu_i}, U_{\mu_j})$. For the mean function, we may simply choose $\mathbb{E}[U_\mu] = \mu$. For the covariance function, we need to ensure that the variance of any combination of random variables in the Gaussian process is non-negative. This property is satisfied by kernel functions. We can therefore define

$$\text{Cov}(U_{\mu_i}, U_{\mu_j}) := k(\mu_i, \mu_j),$$

for some kernel k . Equipped with such a mapping from real numbers to random variables, we need a measure of similarity between random variables. A natural choice is their **correlation**

$$\text{corr}(U_{\mu_i}, U_{\mu_j}) := \frac{\text{Cov}(U_{\mu_i}, U_{\mu_j})}{\sqrt{\text{Var}(U_{\mu_i}) \text{Var}(U_{\mu_j})}} \in [0, 1].$$

We therefore obtain

$$\begin{aligned} \text{corr}(U_{\mu_i}, U_{\mu_j}) &= \frac{k(\mu_1, \mu_2)}{\sqrt{k(\mu_1, \mu_1)k(\mu_2, \mu_2)}} \\ &= \frac{\langle \phi(\mu_1), \phi(\mu_2) \rangle}{\|\phi(\mu_1)\| \|\phi(\mu_2)\|}, \end{aligned}$$

which coincides with the **cosine similarity** measure we saw before. In the particular case $K = 2$ and when $k_\sigma(\mu_1, \mu_2) = \kappa(\mu_1 - \mu_2)$, we then recover the previous heuristically-defined soft equality operator

$$\text{eq}_\sigma(\mu_1, \mu_2) = \text{corr}(U_{\mu_1}, U_{\mu_2}) = \frac{\kappa(\mu_1 - \mu_2)}{\kappa(0)}.$$

5.4 Logical operators

Logical operators can be used to perform Boolean algebra. Formally, we can define them as functions from $\{0, 1\} \times \{0, 1\}$ to $\{0, 1\}$. The **and** (logical conjunction a.k.a. logical product), **or** (logical disjunction a.k.a. logical addition) and **not** (logical negation a.k.a. logical complement) operators, for example, are defined by

$$\begin{aligned} \text{and}(\pi, \pi') &:= \begin{cases} 1 & \text{if } \pi = \pi' = 1 \\ 0 & \text{otherwise} \end{cases} \\ \text{or}(\pi, \pi') &:= \begin{cases} 1 & \text{if } 1 \in \{\pi, \pi'\} \\ 0 & \text{otherwise} \end{cases} \\ \text{not}(\pi) &:= \begin{cases} 0 & \text{if } \pi = 1 \\ 1 & \text{if } \pi = 0 \end{cases}. \end{aligned}$$

Classical properties of these operators include

- Commutativity:

$$\begin{aligned} \text{and}(\pi, \pi') &= \text{and}(\pi', \pi) \\ \text{or}(\pi, \pi') &= \text{or}(\pi', \pi) \end{aligned}$$

- Associativity:

$$\begin{aligned} \text{and}(\pi, \text{and}(\pi', \pi'')) &= \text{and}(\text{and}(\pi, \pi'), \pi'') \\ \text{or}(\pi, \text{or}(\pi', \pi'')) &= \text{or}(\text{or}(\pi, \pi'), \pi'') \end{aligned}$$

- Distributivity of **and** over **or**:

$$\text{and}(\pi, \text{or}(\pi', \pi'')) = \text{or}(\text{and}(\pi, \pi'), \text{and}(\pi, \pi''))$$

- Neutral element:

$$\begin{aligned}\text{and}(\pi, 1) &= \pi \\ \text{or}(\pi, 0) &= \pi\end{aligned}$$

- De Morgan's laws:

$$\begin{aligned}\text{not}(\text{or}(\pi, \pi')) &= \text{and}(\text{not}(\pi), \text{not}(\pi')) \\ \text{not}(\text{and}(\pi, \pi')) &= \text{or}(\text{not}(\pi), \text{not}(\pi')).\end{aligned}$$

More generally, for a binary vector $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K) \in \{0, 1\}^K$, we can define **all** (universal quantification, \forall) and **any** (existential quantification, \exists) operators, which are functions from $\{0, 1\}^K$ to $\{0, 1\}$, as

$$\text{all}(\boldsymbol{\pi}) := \begin{cases} 1 & \text{if } \pi_1 = \dots = \pi_K = 1 \\ 0 & \text{otherwise} \end{cases}$$

and

$$\text{any}(\boldsymbol{\pi}) := \begin{cases} 1 & \text{if } 1 \in \{\pi_1, \dots, \pi_K\} \\ 0 & \text{otherwise} \end{cases}.$$

5.5 Continuous extensions of logical operators

5.5.1 Probabilistic continuous extension

We can equivalently write the **and**, **or** and **not** operators as

$$\begin{aligned}\text{and}(\pi, \pi') &= \pi \cdot \pi' \\ \text{or}(\pi, \pi') &= \pi + \pi' - \pi \cdot \pi' \\ \text{not}(\pi) &= 1 - \pi.\end{aligned}$$

These are **extensions** of the previous definitions: we can use them as functions from $[0, 1] \times [0, 1] \rightarrow [0, 1]$, as illustrated in Fig. 5.3. This means that we can use the soft comparison operators defined in Section 5.2 to obtain $\pi, \pi' \in [0, 1]$. Likewise, we can define continuous extensions of

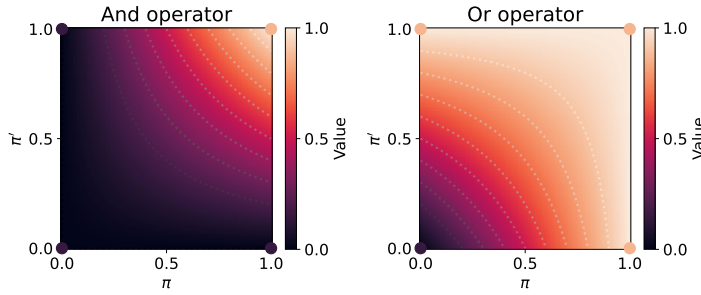


Figure 5.3: The Boolean **and** and **or** operators are functions from $\{0, 1\} \times \{0, 1\}$ to $\{0, 1\}$ (corners in the figure) but their continuous extensions $\text{and}(\pi, \pi') := \pi \cdot \pi'$ as well as $\text{or}(\pi, \pi') := \pi + \pi' - \pi \cdot \pi'$ define a function from $[0, 1] \times [0, 1]$ to $[0, 1]$.

all and **any**, which are functions from $[0, 1]^K$ to $[0, 1]$, as

$$\begin{aligned} \text{all}(\boldsymbol{\pi}) &= \prod_{i=1}^K \pi_i \\ \text{any}(\boldsymbol{\pi}) &= 1 - \prod_{i=1}^K (1 - \pi_i). \end{aligned}$$

From a probabilistic perspective, if we let Y and Y' to be two independent random variables distributed according to **Bernoulli distributions** with parameter π and π' , then

$$\begin{aligned} \text{and}(\pi, \pi') &= \mathbb{P}(Y = 1 \cap Y' = 1) = \mathbb{P}(Y = 1) \cdot \mathbb{P}(Y' = 1) \\ \text{or}(\pi, \pi') &= \mathbb{P}(Y = 1 \cup Y' = 1) \\ &= \mathbb{P}(Y = 1) + \mathbb{P}(Y' = 1) - \mathbb{P}(Y = 1 \cap Y' = 1) \\ &= \mathbb{P}(Y = 1) + \mathbb{P}(Y' = 1) - \mathbb{P}(Y = 1)\mathbb{P}(Y' = 1) \\ \text{not}(\pi) &= \mathbb{P}(Y \neq 1) = 1 - \mathbb{P}(Y = 1). \end{aligned}$$

In probability theory, these correspond to the product rule of two independent variables, the addition rule, and the complement rule.

Likewise, if we let $Y = (Y_1, \dots, Y_K) \in \{0, 1\}^K$ be a random variable distributed according to a **multivariate Bernoulli distribution** with

parameters $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$, then

$$\begin{aligned}
 \text{all}(\boldsymbol{\pi}) &= \mathbb{P}(Y_1 = 1 \cap \dots \cap Y_K = 1) \\
 &= \prod_{i=1}^K \mathbb{P}(Y_i = 1) \\
 \text{any}(\boldsymbol{\pi}) &= \mathbb{P}(Y_1 = 1 \cup \dots \cup Y_K = 1) \\
 &= 1 - \mathbb{P}(\neg(Y_1 = 1 \cup \dots \cup Y_K = 1)) \\
 &= 1 - \mathbb{P}(Y_1 \neq 1 \cap \dots \cap Y_K \neq 1) \\
 &= 1 - \prod_{i=1}^K (1 - \mathbb{P}(Y_i = 1)).
 \end{aligned}$$

These are the chain rule of probability and the addition rule of probability for K independent variables.

5.5.2 Triangular norms and co-norms

More generally, in the **fuzzy logic** literature (Klir and Yuan, 1995; Jayaram and Baczynski, 2008), the concepts of triangular norms and co-norms have been introduced to provide continuous relaxations of the **and** and **or** operators, respectively.

Definition 5.1 (Triangular norms and conorms). A triangular norm, a.k.a. t-norm, is a function from $[0, 1] \times [0, 1]$ to $[0, 1]$ which is commutative, associative, neutral w.r.t. 1 and is monotone, meaning that $t(\pi, \pi') \leq t(\tau, \tau')$ for all $\pi \leq \tau$ and $\pi' \leq \tau'$. A triangular conorm, a.k.a. t-conorm, is defined similarly but is neutral w.r.t. 0.

The previously-defined probabilistic extensions of **and** and **or** are examples of triangle norm and conorm. More examples are given in Table 5.1. Thanks to the associative property of these operators, we can generalize them to vectors $\boldsymbol{\pi} \in [0, 1]^K$ to define continuous extensions of the **all** and **any** operators, as shown in Table 5.2. For more examples and analysis, see for instance van Krieken (2024, Chapters 2 and 3).

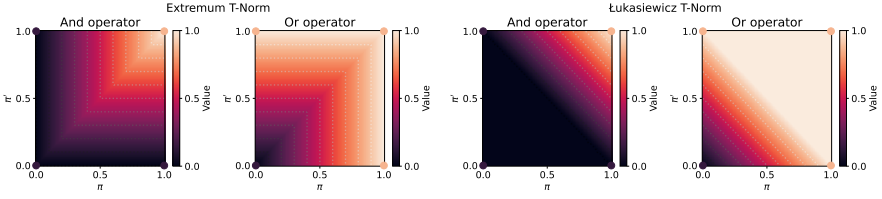


Figure 5.4: Alternative relaxations of the Boolean **and** and **or** operators using triangular norms (t-norms).

Table 5.1: Examples of triangular norms and conorms, which are continuous relaxations of the **and** and **or** operators, respectively. More instances can be obtained by smoothing out the min and max operators.

	t-norm (relaxed and)	t-conorm (relaxed or)
Probabilistic	$\pi \cdot \pi'$	$\pi + \pi' - \pi \cdot \pi'$
Extremum	$\min(\pi, \pi')$	$\max(\pi, \pi')$
Łukasiewicz	$\max(\pi + \pi' - 1, 0)$	$\min(\pi + \pi', 1)$

5.6 If-else statements

An if-else statement executes different code depending on a condition. Formally, we can define the ifelse: $\{0, 1\} \times \mathcal{V} \times \mathcal{V} \rightarrow \mathcal{V}$ function by

$$\text{ifelse}(\pi, \mathbf{v}_1, \mathbf{v}_0) := \begin{cases} \mathbf{v}_1 & \text{if } \pi = 1 \\ \mathbf{v}_0 & \text{if } \pi = 0 \end{cases} \quad (5.3)$$

$$= \pi \cdot \mathbf{v}_1 + (1 - \pi) \cdot \mathbf{v}_0.$$

The π variable is called the **predicate**. It is a **binary** (Boolean) variable, making the function ifelse undefined if $\pi \notin \{0, 1\}$. The function is therefore discontinuous and nondifferentiable w.r.t. $\pi \in \{0, 1\}$. On the other hand, $\mathbf{v}_0 \in \mathcal{V}$ and $\mathbf{v}_1 \in \mathcal{V}$, which correspond to the false and true **branches**, can be **continuous** variables. If $\pi = 1$, the function is linear w.r.t. \mathbf{v}_1 and constant w.r.t. \mathbf{v}_0 . Conversely, if $\pi = 0$, the function is linear w.r.t. \mathbf{v}_0 and constant w.r.t. \mathbf{v}_1 . We now discuss how to differentiate through ifelse.

Table 5.2: Continuous extensions of the **all** and **any** operators.

	All (\forall)	Any (\exists)
Probabilistic	$\prod_{i=1}^K \pi_i$	$1 - \prod_{i=1}^K (1 - \pi_i)$
Extremum	$\min(\pi_1, \dots, \pi_K)$	$\max(\pi_1, \dots, \pi_K)$
Łukasiewicz	$\max(\sum_{i=1}^K \pi_i - (K - 1), 0)$	$\min(\sum_{i=1}^K \pi_i, 1)$

5.6.1 Differentiating through branch variables

For $\pi \in \{0, 1\}$ fixed, $\text{ifelse}(\pi, \mathbf{v}_1, \mathbf{v}_0)$ is a valid function w.r.t. $\mathbf{v}_1 \in \mathcal{V}$ and $\mathbf{v}_0 \in \mathcal{V}$, and can therefore be used as a node in a computational graph (Section 8.3). Due to the linearity w.r.t. \mathbf{v}_1 and \mathbf{v}_0 , we obtain that the Jacobians w.r.t. \mathbf{v}_1 and \mathbf{v}_0 are

$$\begin{aligned} \partial_{\mathbf{v}_0} \text{ifelse}(\pi, \mathbf{v}_1, \mathbf{v}_0) &:= \begin{cases} 0 & \text{if } \pi = 1 \\ I & \text{if } \pi = 0 \end{cases} \\ &= (1 - \pi) \cdot I \end{aligned}$$

and

$$\begin{aligned} \partial_{\mathbf{v}_1} \text{ifelse}(\pi, \mathbf{v}_1, \mathbf{v}_0) &:= \begin{cases} I & \text{if } \pi = 1 \\ 0 & \text{if } \pi = 0 \end{cases} \\ &= \pi \cdot I, \end{aligned}$$

where I is the identity matrix of appropriate size. Most of the time, if-else statements are composed with other functions. Let $g_1: \mathcal{U}_1 \rightarrow \mathcal{V}$ and $g_0: \mathcal{U}_0 \rightarrow \mathcal{V}$ be differentiable functions. We then define $\mathbf{v}_1 := g_1(\mathbf{u}_1)$ and $\mathbf{v}_0 := g_0(\mathbf{u}_0)$, where $\mathbf{u}_1 \in \mathcal{U}_1$ and $\mathbf{u}_0 \in \mathcal{U}_0$. The composition of ifelse , g_1 and g_0 is then the function $f: \{0, 1\} \times \mathcal{U}_1 \times \mathcal{U}_0 \rightarrow \mathcal{V}$ defined by

$$\begin{aligned} f(\pi, \mathbf{u}_1, \mathbf{u}_0) &:= \text{ifelse}(\pi, g_1(\mathbf{u}_1), g_0(\mathbf{u}_0)) \\ &= \pi \cdot g_1(\mathbf{u}_1) + (1 - \pi) \cdot g_0(\mathbf{u}_0). \end{aligned}$$

We obtain that the Jacobians are

$$\partial_{\mathbf{u}_1} f(\pi, \mathbf{u}_1, \mathbf{u}_0) = \pi \cdot \partial g_1(\mathbf{u}_1)$$

and

$$\partial_{\mathbf{u}_0} f(\pi, \mathbf{u}_1, \mathbf{u}_0) = (1 - \pi) \partial g_0(\mathbf{u}_0).$$

As long as g_1 and g_0 are differentiable functions, we can therefore differentiate through the branch variables \mathbf{u}_1 and \mathbf{u}_0 without any issue. More problematic is the predicate variable π , as we now discuss.

5.6.2 Differentiating through predicate variables

The predicate variable π is binary and therefore cannot be differentiated directly. However, π can be the output of a comparison operator. For example, suppose we want to express the function $f_h: \mathbb{R} \times \mathcal{U}_1 \times \mathcal{U}_0 \rightarrow \mathcal{V}$ defined by

$$f_h(p, \mathbf{u}_1, \mathbf{u}_0) := \begin{cases} g_1(\mathbf{u}_1) & \text{if } p \geq 0 \\ g_0(\mathbf{u}_0) & \text{otherwise} \end{cases}.$$

Using our notation, this can be rewritten as

$$\begin{aligned} f_h(p, \mathbf{u}_1, \mathbf{u}_0) &:= \text{ifelse}(\text{gt}(p, 0), g_1(\mathbf{u}_1), g_0(\mathbf{u}_0)) \\ &= \text{ifelse}(\text{step}(p), g_1(\mathbf{u}_1), g_0(\mathbf{u}_0)) \\ &= \text{step}(p)g_1(\mathbf{u}_1) + (1 - \text{step}(p))g_0(\mathbf{u}_0). \end{aligned}$$

The Heaviside step function has a **discontinuity** at $p = 0$, but it is continuous and differentiable with derivative $\text{step}'(p) = 0$ for all $p \neq 0$. The function f_h therefore has **null derivative** w.r.t. $p \neq 0$,

$$\begin{aligned} \partial_p f_h(p, \mathbf{u}_1, \mathbf{u}_0) &= \partial_1 f_h(p, \mathbf{u}_1, \mathbf{u}_0) \\ &= \text{step}'(p)(g_1(\mathbf{u}_1) - g_0(\mathbf{u}_0)) \\ &= \mathbf{0}. \end{aligned}$$

In other words, while f_h has **well-defined** derivatives w.r.t. p for $p \neq 0$, the derivatives are **uninformative**. As another example, let us now consider the function

$$g_h(\mathbf{u}_1, \mathbf{u}_0) := f_h(t(\mathbf{u}_1), \mathbf{u}_1, \mathbf{u}_0),$$

for some differentiable function t . This time, \mathbf{u}_1 influences both the predicate and the true branch. Then, using Proposition 2.8, we obtain

$$\begin{aligned} \partial_{\mathbf{u}_1} g_h(\mathbf{u}_1, \mathbf{u}_0) &= \partial t(\mathbf{u}_1) \partial_1 f_h(t(\mathbf{u}_1), \mathbf{u}_1, \mathbf{u}_0) + \partial_2 f_h(t(\mathbf{u}_1), \mathbf{u}_1, \mathbf{u}_0) \\ &= \partial_2 f_h(t(\mathbf{u}_1), \mathbf{u}_1, \mathbf{u}_0). \end{aligned}$$

In other words, the derivatives of the predicate $t(\mathbf{u}_1)$ do not influence the derivatives of g_h .

5.6.3 Continuous relaxations

Fortunately, we recall that

$$\text{ifelse}(\pi, \mathbf{v}_1, \mathbf{v}_0) = \pi \cdot \mathbf{v}_1 + (1 - \pi) \cdot \mathbf{v}_0.$$

This function is perfectly well-defined, even if $\pi \in [0, 1]$, instead of $\pi \in \{0, 1\}$. That is, this definition is an **extension** of Eq. (5.3) from the discrete set $\{0, 1\}$ to the continuous unit segment $[0, 1]$. We saw that

$$\text{gt}(a, b) \approx \text{gt}_\sigma(a, b) := \text{sigmoid}_\sigma(a - b) \in [0, 1],$$

where we use sigmoid_σ to denote a differentiable S-shaped function mapping \mathbb{R} to $[0, 1]$. For instance, we can use the logistic function or the standard Gaussian's CDF. If we now define

$$\begin{aligned} f_s(p, \mathbf{u}_1, \mathbf{u}_0) &:= \text{ifelse}(\text{gt}_\sigma(p), g_1(\mathbf{u}_1), g_0(\mathbf{u}_0)) \\ &= \text{ifelse}(\text{sigmoid}_\sigma(p), g_1(\mathbf{u}_1), g_0(\mathbf{u}_0)) \\ &= \text{sigmoid}_\sigma(p)g_1(\mathbf{u}_1) + (1 - \text{sigmoid}_\sigma(p))g_0(\mathbf{u}_0), \end{aligned} \quad (5.4)$$

the Jacobian becomes

$$\partial_p f_s(p, \mathbf{u}_1, \mathbf{u}_0) = \text{sigmoid}'_\sigma(p)(g_1(\mathbf{u}_1) - g_0(\mathbf{u}_0)).$$

If $\text{sigmoid}_\sigma = \text{logistic}(\cdot/\sigma)$ or $\text{sigmoid}_\sigma = \Phi(\cdot/\sigma)$, the Jacobian is **non-null everywhere**, allowing gradients to backpropagate through the computational graph. This is an example of smoothing as studied in Part IV.

Probabilistic perspective

From a probabilistic perspective, we can view Eq. (5.4) as the expectation of $g_i(\mathbf{u}_i)$, where $i \in \{0, 1\}$ is a binary random variable distributed according to a **Bernoulli distribution** with parameter $\pi = \text{sigmoid}_\sigma(p)$:

$$f_s(p, \mathbf{u}_1, \mathbf{u}_0) = \mathbb{E}_{i \sim \text{Bernoulli}(\text{sigmoid}_\sigma(p))} [g_i(\mathbf{u}_i)].$$

Taking the expectation over the two possible branches makes the function differentiable with respect to p , since $\text{sigmoid}_\sigma(p)$ is differentiable. Of course, this comes at the cost of evaluating both branches, instead

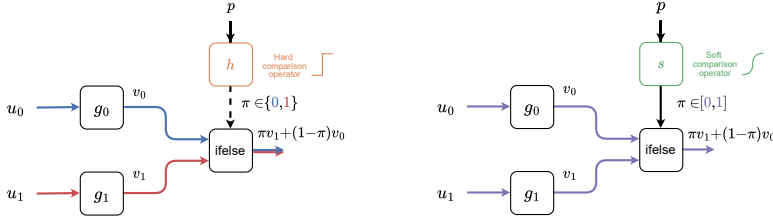


Figure 5.5: Computation graphs of programs using if-else statements with either hard or soft comparison operators. By using a hard comparison operator (step function, left panel) the predicate π is a discrete variable (represented by a dashed line). Depending on the value (0 or 1) of the predicate π , only one branch (red or blue) contributes to the output. Derivatives along a path of continuous variables (dense lines) can be computed. However, discrete variables such as the predicate prevent the propagation of meaningful derivatives. By using a soft comparison operator (sigmoid, right panel), the predicate is a continuous variable and derivatives with respect to the input p can be taken. In this case both branches (corresponding to g_0 and g_1) contribute to the output and therefore need to be evaluated.

of a single one. The probabilistic perspective suggests that we can also compute the variance if needed as

$$\begin{aligned} & \mathbb{V}_{i \sim \text{Bernoulli}(\text{sigmoid}_\sigma(p))} [g_i(\mathbf{u}_i)] \\ &= \mathbb{E}_{i \sim \text{Bernoulli}(\text{sigmoid}_\sigma(p))} \left[(f_s(p, \mathbf{u}_1, \mathbf{u}_0) - g_i(\mathbf{u}_i))^2 \right]. \end{aligned}$$

The probabilistic viewpoint also suggests different scales at which a smoothing can be defined as illustrated in Fig. 5.6.

Another perspective (Petersen *et al.*, 2021) is based on the **logistic distribution**. Indeed, if P is a random variable following a logistic distribution with mean p and scale 1, we saw in Remark 3.1 that the CDF is $\mathbb{P}(P \leq 0) = \text{logistic}(-p) = 1 - \text{logistic}(p)$ and therefore

$$\begin{aligned} f_s(p, \mathbf{u}_1, \mathbf{u}_0) &= \text{ifelse}(\text{logistic}(p), g_1(\mathbf{u}_1), g_0(\mathbf{u}_0)) \\ &= \text{logistic}(p)g_1(\mathbf{u}_1) + (1 - \text{logistic}(p))g_0(\mathbf{u}_0) \\ &= \mathbb{P}(P > 0) \cdot g_1(\mathbf{u}_1) + \mathbb{P}(P \leq 0) \cdot g_0(\mathbf{u}_0). \end{aligned}$$

Remark 5.1 (Global versus local smoothing). Consider the function

$$f(x, y, z) := \begin{cases} y & \text{if } a \leq x \leq b \\ z & \text{otherwise} \end{cases}.$$

The derivatives w.r.t. y and z are well-defined. The derivative w.r.t. x on the other hand is not well-defined since it involves comparison operators and the logical operator and. Using our notation, we can rewrite the function as

$$f(x, y, z) = \text{ifelse}(\text{and}(\text{gt}(x, a), \text{lt}(x, b)), y, z).$$

A local smoothing approach consists in replacing gt and lt by gt_σ and lt_σ locally in the program:

$$\begin{aligned} f_\sigma^{\text{loc}}(x, y, z) &:= \text{ifelse}(\text{and}(\text{gt}_\sigma(x, a), \text{lt}_\sigma(x, b)), y, z) \\ &= \pi_a \pi_b y + (1 - \pi_a \pi_b) z \end{aligned}$$

where

$$\begin{aligned} \pi_a &:= \text{sigmoid}_\sigma(x - a) \\ \pi_b &:= \text{sigmoid}_\sigma(b - x), \end{aligned}$$

for any sigmoid function sigmoid_σ . A global smoothing approach instead uses the expectation of the entire program

$$\begin{aligned} f_\sigma^{\text{glob}}(x, y, z) &:= \mathbb{E}_Z[\text{ifelse}(\text{and}(\text{gt}(x + \sigma Z, a), \text{lt}(x + \sigma Z, b)), y, z)] \\ &= \text{ifelse}(\pi, y, z) \end{aligned}$$

where

$$\begin{aligned} \pi &:= \mathbb{E}_Z[\text{and}(\text{gt}(x + \sigma Z, a), \text{lt}(x + \sigma Z, b))] \\ &= \mathbb{P}(a \leq x + \sigma Z \leq b) \\ &= \text{sigmoid}_\sigma(b - x) - \text{sigmoid}_\sigma(a - x) \\ &= \pi_b - \pi_a, \end{aligned}$$

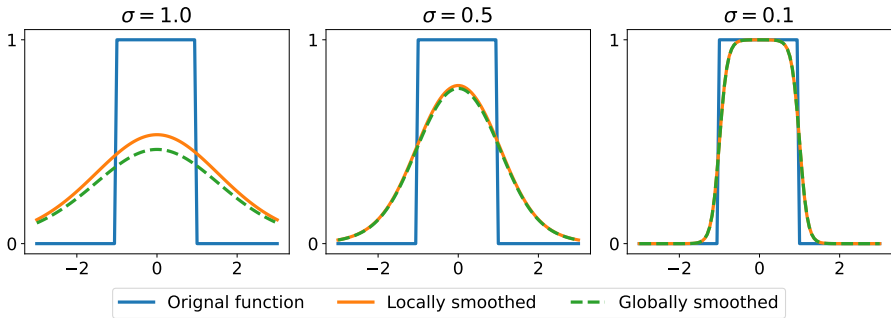


Figure 5.6: Global versus local smoothing approaches on a gate function $f(x) := 1$ if $x \in [-1, 1]$, and $f(x) := 0$ otherwise. In our notation, we can write $f(x) = \text{ifelse}(\text{and}(\text{gt}(x, -1), \text{lt}(x, 1)), 1, 0)$. A local approach smooths out gt and lt separately. A global approach uses the expectation of the whole program, see Remark 5.1. We observe that, though the approaches differ for large σ , they quickly coincide for smaller σ .

for sigmoid_σ the CDF of σZ . We therefore obtain

$$f_\sigma^{\text{glob}}(x, y, z) = (\pi_b - \pi_a)y + (1 - (\pi_b - \pi_a))z.$$

The difference stems from the fact that the local approach smooths out $a \leq x$ and $x \leq b$ independently (treating $\mathbf{1}_{X \geq a}$ and $\mathbf{1}_{X \leq b}$ as independent random variables), while the global approach smooths out $a \leq x \leq b$ simultaneously. In practice, both approaches approximate the original function well as $\sigma \rightarrow 0$ and coincide for σ sufficiently small as illustrated in Fig. 5.6.

5.7 Else-if statements

In the previous section, we focused on if-else statements: conditionals with only two branches. We now generalize our study to conditionals including else-if statements, that have K branches.

5.7.1 Encoding K branches

For conditionals with only 2 branches, we encoded the branch that the conditional needs to take using the binary variable $\pi \in \{0, 1\}$. For

conditionals with K branches, we need a way to encode which of the K branches the conditional needs to take. To do so, we can use a vector $\boldsymbol{\pi} \in \{\mathbf{e}_1, \dots, \mathbf{e}_K\}$, where \mathbf{e}_i denotes the standard basis vector (a.k.a. one-hot vector)

$$\mathbf{e}_i := (0, \dots, \underbrace{1}_i, \dots, 0),$$

a vector with a single one in the coordinate i and $K - 1$ zeros. The vector \mathbf{e}_i is the encoding of a **categorical variable** $i \in [K]$.

Combining booleans

To form, such a vector $\boldsymbol{\pi} \in \{\mathbf{e}_1, \dots, \mathbf{e}_K\}$, we can combine the previously-defined comparison and logical operators to define $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$. However, we need to ensure that only one π_i is non-zero. We give an example in Example 5.1.

Argmax and argmin operators

Another way to form $\boldsymbol{\pi}$ is to use the **argmax** and **argmin** operators

$$\begin{aligned} \operatorname{argmax}(\mathbf{p}) &:= \arg \max_{\boldsymbol{\pi} \in \{\mathbf{e}_1, \dots, \mathbf{e}_K\}} \langle \boldsymbol{\pi}, \mathbf{p} \rangle \\ \operatorname{argmin}(\mathbf{p}) &:= \arg \min_{\boldsymbol{\pi} \in \{\mathbf{e}_1, \dots, \mathbf{e}_K\}} \langle \boldsymbol{\pi}, \mathbf{p} \rangle = \operatorname{argmax}(-\mathbf{p}). \end{aligned}$$

They can be seen as a natural generalization of the greater than and less than operators. In case of ties, we break them arbitrarily.

5.7.2 Conditionals

We can now express a conditional statement as the function $\operatorname{cond}: \{\mathbf{e}_1, \dots, \mathbf{e}_K\} \times \mathcal{V}^K \rightarrow \mathcal{V}$ defined by

$$\begin{aligned} \operatorname{cond}(\boldsymbol{\pi}, \mathbf{v}_1, \dots, \mathbf{v}_K) &:= \begin{cases} \mathbf{v}_1 & \text{if } \boldsymbol{\pi} = \mathbf{e}_1 \\ \vdots & \\ \mathbf{v}_K & \text{if } \boldsymbol{\pi} = \mathbf{e}_K \end{cases} \\ &= \sum_{i=1}^K \pi_i \mathbf{v}_i. \end{aligned} \tag{5.5}$$

Similarly as for the ifelse function, the cond function is discontinuous and nondifferentiable w.r.t. $\pi \in \{e_1, \dots, e_K\}$. However, given $\pi = e_i$ fixed for some i , the function is linear in v_i and constant in v_j for $j \neq i$. We illustrate how to express a simple example, using this formalism.

Example 5.1 (Soft-thresholding operator). The soft-thresholding operator (see also Section 16.4) is a commonly-used operator to promote sparsity. It is defined by

$$\text{SoftThreshold}(u, \lambda) := \begin{cases} 0 & \text{if } |u| \leq \lambda \\ u - \lambda & \text{if } u \geq \lambda \\ u + \lambda & \text{if } u \leq -\lambda \end{cases}.$$

To express it in our formalism, we can define $\pi \in \{e_1, e_2, e_3\}$ using comparison operators as

$$\begin{aligned} \pi &:= (\text{lt}(|u|, \lambda), \text{gt}(u, \lambda), \text{lt}(u, -\lambda)) \\ &= (\text{step}(\lambda - |u|), \text{step}(u - \lambda), \text{step}(-u - \lambda)). \end{aligned}$$

Equivalently, we can also define π using an argmax operator as

$$\pi := \text{argmax}((\lambda - |u|, u - \lambda, -u - \lambda)).$$

In case of ties, which happens at $|u| = \lambda$, we keep only one non-zero coordinate in π . We can then rewrite the operator as

$$\text{SoftThreshold}(u, \lambda) = \text{cond}(\pi, 0, u - \lambda, u + \lambda).$$

As we will see, replacing argmax with softargmax induces a categorical distribution over the three possible branches. The mean value can be seen as a smoothed out version of the operator, and we can also compute the standard deviation, as illustrated in Fig. 5.7.

5.7.3 Differentiating through branch variables

For π fixed, $\text{cond}(\pi, v_1, \dots, v_K)$ is a valid function w.r.t. v_i , and can therefore again be used as a node in a computational graph. Due to the

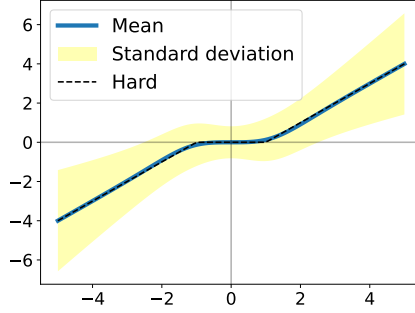


Figure 5.7: A conditional with three branches: the soft-thresholding operator (see Example 5.1). It is a piecewise linear function (dotted black line). Using a softargmax, we can induce a categorical probability distribution over the three branches. The expected value (blue line) can be seen as a smoothed out version of the operator. The induced distribution allows us to also compute the standard deviation.

linearity w.r.t. \mathbf{v}_i , we obtain that the Jacobian w.r.t. \mathbf{v}_i is

$$\partial_{\mathbf{v}_i} \text{cond}(\boldsymbol{\pi}, \mathbf{v}_1, \dots, \mathbf{v}_K) := \begin{cases} I & \text{if } \boldsymbol{\pi} = \mathbf{e}_i \\ \mathbf{0} & \text{if } \boldsymbol{\pi} \neq \mathbf{e}_i \end{cases}.$$

Let $g_i: \mathcal{U}_i \rightarrow \mathcal{V}$ be a differentiable function and $\mathbf{u}_i \in \mathcal{U}_i$. If we define the composition

$$f(\boldsymbol{\pi}, \mathbf{u}_1, \dots, \mathbf{u}_K) := \text{cond}(\boldsymbol{\pi}, g_1(\mathbf{u}_1), \dots, g_K(\mathbf{u}_K)),$$

we then obtain that the Jacobian w.r.t. \mathbf{u}_i is

$$\partial_{\mathbf{u}_i} f(\boldsymbol{\pi}, \mathbf{u}_1, \dots, \mathbf{u}_K) := \begin{cases} \partial g_i(\mathbf{u}_i) & \text{if } \boldsymbol{\pi} = \mathbf{e}_i \\ \mathbf{0} & \text{if } \boldsymbol{\pi} \neq \mathbf{e}_i \end{cases}.$$

As long as the g_i functions are differentiable, we can therefore differentiate through the branch variables \mathbf{u}_i for $\boldsymbol{\pi}$ fixed.

5.7.4 Differentiating through predicate variables

As we saw, $\boldsymbol{\pi}$ can be obtained by combining comparison and logical operators, or it can be obtained by argmax and argmin operators.

We illustrate here why these operators are problematic. For example, suppose we want to express the function

$$f_a(\mathbf{p}, \mathbf{u}_1, \dots, \mathbf{u}_K) := \begin{cases} \mathbf{v}_1 & \text{if } \mathbf{p} = \mathbf{e}_1 \\ \vdots & \\ \mathbf{v}_K & \text{if } \mathbf{p} = \mathbf{e}_K \end{cases}.$$

In our notation, this can be expressed as

$$f_a(\mathbf{p}, \mathbf{u}_1, \dots, \mathbf{u}_K) := \text{cond}(\text{argmax}(\mathbf{p}), g_1(\mathbf{u}_1), \dots, g_K(\mathbf{u}_K)),$$

As for the ifelse case, the Jacobian w.r.t. \mathbf{p} is null almost everywhere,

$$\partial_{\mathbf{p}} f_a(\mathbf{p}, \mathbf{u}_1, \dots, \mathbf{u}_K) = \mathbf{0}.$$

5.7.5 Continuous relaxations

Similarly to the Heaviside step function, the argmax and argmin functions are piecewise constant, with discontinuities in case of ties. Their Jacobian are zero almost everywhere, and undefined in case of ties. Therefore, while their Jacobian is well-defined almost everywhere, they are uninformative and prevent gradient backpropagation. We can replace the argmax with a **softargmax**

$$\text{softargmax}(\mathbf{p}) := \frac{\exp(\mathbf{p})}{\sum_{i=1}^K \exp(p_i)} \in \Delta^K$$

and similarly

$$\text{softargmin}(\mathbf{p}) := \text{softargmax}(-\mathbf{p}) \in \Delta^K.$$

Other relaxations of the argmax are possible, as discussed in Section 13.7. See also Section 14.5.3 for the perturbation perspective.

Fortunately, the definition

$$\text{cond}(\boldsymbol{\pi}, \mathbf{v}_1, \dots, \mathbf{v}_K) = \sum_{i=1}^K \pi_i \mathbf{v}_i$$

is perfectly valid if we use $\boldsymbol{\pi} \in \Delta^K$ instead of $\boldsymbol{\pi} \in \{\mathbf{e}_1, \dots, \mathbf{e}_K\}$, and

can therefore be seen as an **extension** of Eq. (5.5). If we now define

$$\begin{aligned} f_s(\mathbf{p}, \mathbf{u}_1, \dots, \mathbf{u}_K) &:= \text{cond}(\text{softargmax}(\mathbf{p}), g_1(\mathbf{u}_1), \dots, g_K(\mathbf{u}_K)) \\ &= \sum_{i=1}^K [\text{softargmax}(\mathbf{p})]_i \cdot g_i(\mathbf{u}_i), \end{aligned} \quad (5.6)$$

the Jacobian becomes

$$\partial_p f_s(\mathbf{p}, \mathbf{u}_1, \dots, \mathbf{u}_K) = \partial \text{softargmax}(\mathbf{p})(g_1(\mathbf{u}_1), \dots, g_K(\mathbf{u}_K)),$$

which is **non-null everywhere**, allowing gradients to backpropagate through the computational graph.

Probabilistic perspective

From a probabilistic perspective, we can view Eq. (5.6) as the expectation of $g_i(\mathbf{u}_i)$, where $i \in [K]$ is a categorical random variable distributed according to a **categorical distribution** with parameter $\pi = \text{softargmax}(\mathbf{p})$:

$$f_s(\mathbf{p}, \mathbf{u}_1, \dots, \mathbf{u}_K) = \mathbb{E}_{i \sim \text{Categorical}(\text{softargmax}(\mathbf{p}))} [g_i(\mathbf{u}_i)].$$

Taking the expectation over the K possible branches makes the function differentiable with respect to \mathbf{p} , at the cost of evaluating all branches, instead of a single one. Similarly as for the if-else case, we can compute the variance if needed as

$$\begin{aligned} &\mathbb{V}_{i \sim \text{Categorical}(\text{softargmax}(\mathbf{p}))} [g_i(\mathbf{u}_i)] \\ &= \mathbb{E}_{i \sim \text{Categorical}(\text{softargmax}(\mathbf{p}))} \left[(f_s(\mathbf{p}, \mathbf{u}_1, \dots, \mathbf{u}_K) - g_i(\mathbf{u}_i))^2 \right]. \end{aligned}$$

This is illustrated in Fig. 5.7.

5.8 For loops

For loops are a control flow for sequentially calling a fixed number K of functions, reusing the output from the previous iteration. In full generality, a for loop can be written as follows.

Algorithm 5.1 $r = \text{forloop}(s_0)$

for $k := 1, \dots, K$ **do**

$s_k := f_k(s_{k-1})$

$r := s_K$

As illustrated in Fig. 5.8, this defines a computation chain. Assuming the functions f_k are all differentiable, this defines a valid computation graph, we can therefore use automatic differentiation to differentiate `forloop` w.r.t. its input s_0 . Feedforward networks, reviewed in Section 4.2, can be seen as **parameterized for loops**, i.e.,

$$f_k(s_{k-1}) := g_k(s_{k-1}, w_k),$$

for some differentiable function g_k .

Example 5.2 (Unrolled gradient descent). Suppose we want to minimize w.r.t. w the function

$$L(w, \lambda) := \frac{1}{N} \sum_{i=1}^N \ell(h(x_i, w), y_i) + \frac{\lambda}{2} \|w\|_2^2.$$

Given an initialization w_0 , gradient descent (Section 16.1) performs iterations of the form

$$w_k = f(w_{k-1}, \gamma_k, \lambda) := w_{k-1} - \gamma_k \nabla_1 L(w_{k-1}, \lambda).$$

Gradient descent can therefore be expressed as a for loop with

$$f_k(w_{k-1}) := f(w_{k-1}, \gamma_k, \lambda).$$

This means that we can differentiate through the iterations of gradient descent, as long as f is differentiable, meaning that L is twice differentiable. This is useful for instance to perform gradient-based optimization of the hyperparameters γ_k or λ . This a special case of bilevel optimization; see also Chapter 11.

Example 5.3 (Bubble sort). Bubble sort is a simple sorting algorithm that works by repeatedly swapping elements if necessary.

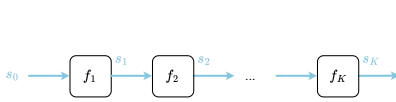


Figure 5.8: A **for loop** forms a computation chain. A feed forward network can be seen as a parameterized for loop, where each function f_k depends on some parameters w_k .

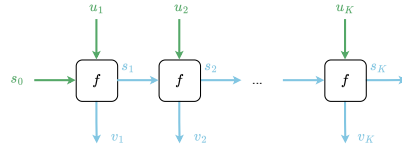


Figure 5.9: Computation graph of the **scan function**. Sequence-to-sequence RNNs can be seen as a parameterized scan function.

Mathematically, swapping two elements i and j can be written as a function from $\mathbb{R}^N \times [N] \times [N]$ to \mathbb{R}^N defined by

$$\text{swap}(\mathbf{v}, i, j) := \mathbf{v} + (v_j - v_i)\mathbf{e}_i + (v_i - v_j)\mathbf{e}_j.$$

We can then write bubble sort as

```

for  $i := 1, \dots, N$  do
  for  $j := 1, \dots, N - i - 1$  do
     $\mathbf{v}' := \text{swap}(\mathbf{v}, j, j + 1)$ 
     $\pi := \text{step}(v_j - v_{j+1})$ 
     $\mathbf{v} \leftarrow \text{ifelse}(\pi, \mathbf{v}', \mathbf{v})$ 

```

Replacing the Heaviside step function with the logistic function gives a smoothed version of the algorithm.

5.9 Scan functions

Scan is a higher-order function (meaning a function of a function) originating from functional programming. It is useful to perform an operation f on individual elements u_k while carrying the result s_k of that operation to the next iteration.

Algorithm 5.2 $r = \text{scan}(s_0, u_1, \dots, u_K)$

```

for  $k := 1, \dots, K$  do
     $s_k, v_k := f(s_{k-1}, u_k)$ 
 $r := (s_K, v_1, \dots, v_K)$ 

```

As illustrated in Fig. 5.9, this again defines a valid computational graph and can be differentiated through using autodiff, assuming the function f is differentiable. Sequence-to-sequence RNNs, reviewed in Section 4.7, can be seen as a **parameterized scan**. An advantage of this abstraction is that parallel scan algorithms have been studied extensively in computer science (Blelloch, 1989; Sengupta *et al.*, 2010).

Example 5.4 (Prefix sum). Scan can be seen as a generalization of the **prefix sum** (a.k.a. **cumulated sum**) from the addition to any binary operation. Indeed, a prefix sum amounts to perform

$$\begin{aligned}
 v_1 &:= u_1 \\
 v_2 &:= u_1 + u_2 \\
 v_3 &:= u_1 + u_2 + u_3 \\
 &\vdots
 \end{aligned}$$

which can be expressed as a scan by defining

$$\begin{aligned}
 v_k &:= s_{k-1} + u_k \\
 f(s_{k-1}, u_k) &:= (v_k, v_k)
 \end{aligned}$$

starting from $s_0 = 0$ (s_K and v_K are redundant in this case).

5.10 While loops

5.10.1 While loops as cyclic graphs

A while loop is a control flow used to repeatedly perform an operation, reusing the output of the previous iteration, until a certain condition is met. Suppose $f: \mathcal{S} \rightarrow \{0, 1\}$ is a function to determine whether to stop ($\pi = 1$) or continue ($\pi = 0$) and $g: \mathcal{S} \rightarrow \mathcal{S}$ is a function for performing

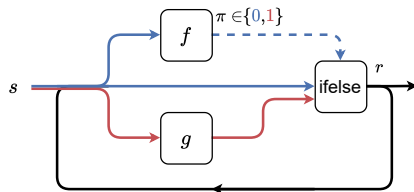


Figure 5.10: A while loop can be represented as a cyclic graph. The while loop stops if $\pi = 1$ and performs another iteration $s \leftarrow g(s)$, $\pi \leftarrow f(s)$ if $\pi = 0$.

an operation. Then, without loss of generality, a while loop can be written as follows.

Algorithm 5.3 $r = \text{whileloop}(s)$

$\pi \leftarrow f(s)$

while $\pi = 0$ **do**

$s \leftarrow g(s)$

$\pi \leftarrow f(s)$

$r := s$

This definition is somewhat cyclic, as we used the **while** keyword. However, we can equivalently rewrite the algorithm recursively.

Algorithm 5.4 $r = \text{whileloop}(s)$

$\pi := f(s)$

if $\pi = 0$ **then**

$r := s$

else

$r := \text{whileloop}(g(s))$

Unlike for loops and scan, the number of iterations of while loops is not known ahead of time, and may even be infinite. In this respect, a while loop can be seen as a **cyclic graph**, as illustrated in Fig. 5.10.

Importance of lazy evaluation

We can also implement Algorithm 5.4 in terms of the ifelse function defined in Section 5.6 as

$$\begin{aligned} r &:= \text{ifelse}(f(s), s, \text{whileloop}(g(s))) \\ &= f(s) \cdot s + (1 - f(s)) \cdot \text{whileloop}(g(s)). \end{aligned}$$

However, to avoid an infinite recursion, it is crucial that ifelse supports **lazy evaluation**. That is, $\text{whileloop}(g(s))$ in the definition above should be evaluated if and only if $\pi = f(s) = 0$. In other words, the fact that $f(s) \in \{0, 1\}$ is crucial to ensure that the recursion is well-defined.

5.10.2 Unrolled while loops

To avoid the issues with unbounded while loops, we can enforce that a while loop stops after T iterations, i.e., we can truncate the while loop. Unrolling Algorithm 5.4 gives (here with $T = 3$)

```

 $\pi_0 := f(s_0)$ 
if  $\pi_0 = 1$  then
   $r := s_0$ 
else
   $s_1 := g(s_0), \pi_1 := f(s_1)$ 
  if  $\pi_1 = 1$  then
     $r := s_1$ 
  else
     $s_2 := g(s_1), \pi_2 := f(s_2)$ 
    if  $\pi_2 = 1$  then
       $r := s_2$ 
    else
       $r := s_3 := g(s_2)$ 

```

Using the ifelse function, we can rewrite it as

$$\begin{aligned} \mathbf{r} = & \text{ifelse}(\pi_0, \\ & \mathbf{s}_0, \\ & \text{ifelse}(\pi_1, \\ & \mathbf{s}_1, \\ & \text{ifelse}(\pi_2, \\ & \mathbf{s}_2, \\ & \mathbf{s}_3))) \end{aligned}$$

which is itself equivalent to

$$\begin{aligned} \mathbf{r} &= \pi_0 \mathbf{s}_0 + (1 - \pi_0) [\pi_1 \mathbf{s}_1 + (1 - \pi_1) [\pi_2 \mathbf{s}_2 + (1 - \pi_2) \mathbf{s}_3]] \\ &= \pi_0 \mathbf{s}_0 + (1 - \pi_0) \pi_1 \mathbf{s}_1 + (1 - \pi_0)(1 - \pi_1) \pi_2 \mathbf{s}_2 + (1 - \pi_0)(1 - \pi_1)(1 - \pi_2) \mathbf{s}_3. \end{aligned}$$

More generally, for $T \in \mathbb{N}$, the formula is

$$\begin{aligned} \mathbf{r} &= \sum_{i=0}^T ((1 - \pi_0) \dots (1 - \pi_{i-1})) \pi_i \mathbf{s}_i \\ &= \sum_{i=0}^T \left(\prod_{j=0}^{i-1} (1 - \pi_j) \right) \pi_i \mathbf{s}_i, \end{aligned}$$

where we defined

$$\mathbf{s}_i := g(\mathbf{s}_{i-1}) := g^i(\mathbf{s}_0) := \underbrace{g \circ \dots \circ g}_{i \text{ times}}(\mathbf{s}_0) \in \mathcal{S}$$

$$\pi_i := f(\mathbf{s}_i) \in \{0, 1\}.$$

See also (Petersen *et al.*, 2021). If we further define the shorthand notation

$$\begin{aligned} \tilde{\pi}_0 &:= \pi_0 \\ \tilde{\pi}_i &:= \left(\prod_{j=0}^{i-1} 1 - \pi_j \right) \pi_i \quad i \in \{1, \dots, T\}, \end{aligned}$$

so that $\tilde{\boldsymbol{\pi}} := (\tilde{\pi}_0, \tilde{\pi}_1, \dots, \tilde{\pi}_T) \in \Delta^{T+1}$ is a discrete probability distribution containing the probabilities to stop at each of the T iterations, we

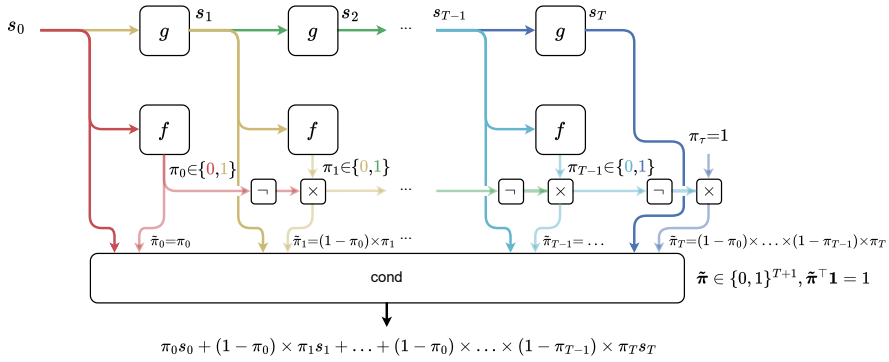


Figure 5.11: Computation graph of an unrolled truncated while loop. As in Fig. 5.5, we depict continuous variables in dense lines and discrete variables in dashed lines. The output of a while loop with at most T iterations can be written as a conditional with $T + 1$ branches, $\text{cond}(\tilde{\pi}, s_0, \dots, s_T) = \sum_{t=0}^T \tilde{\pi}_t s_t$.

can rewrite the output of a truncated while using a conditional,

$$\mathbf{r} = \text{cond}(\tilde{\pi}, s_0, s_1, \dots, s_T) = \sum_{t=0}^T \tilde{\pi}_t s_t.$$

This is illustrated in Fig. 5.11.

Example 5.5 (Computing the square root using Newton's method).

Computing the square root \sqrt{x} of a real number $x > 0$ can be cast as a root finding problem, which we can solve using Newton's method. Starting from an initialization s_0 , the iterations read

$$s_{i+1} := g(s_i) := \frac{1}{2} \left(s_i + \frac{x}{s_i} \right).$$

To measure the error on iteration i , we can define

$$\varepsilon(s_i) := \frac{1}{2}(s_i^2 - x)^2.$$

As a stopping criterion, we can then use

$$\begin{aligned}\pi_i &:= \begin{cases} 1 & \text{if } \varepsilon(s_i) \leq \tau \\ 0 & \text{otherwise} \end{cases} \\ &= \text{step}(\tau - \varepsilon(s_i)),\end{aligned}$$

where $0 < \tau \ll 1$ is an error tolerance and step is the Heaviside step function.

5.10.3 Markov chain perspective

Given the function $g: \mathcal{S} \rightarrow \mathcal{S}$ and the initialization $\mathbf{s}_0 \in \mathcal{S}$, a while loop can only go through a discrete set of values $\mathbf{s}_0, \mathbf{s}_1, \mathbf{s}_2, \dots$ defined by $\mathbf{s}_i = g(\mathbf{s}_{i-1})$. This set is potentially countably infinite if the while loop is unbounded, and finite if the while loop is guaranteed to stop. Whether the loop moves from the state \mathbf{s}_i to the state \mathbf{s}_{i+1} , or stays at \mathbf{s}_i , is determined by the stopping criterion $\pi_i \in \{0, 1\}$. To model the state of the while loop, we can then consider a **Markov chain** with a discrete space $\{\mathbf{s}_0, \mathbf{s}_1, \mathbf{s}_2, \dots\}$, which we can always identify with $\{0, 1, 2, \dots\}$, with transition probabilities

$$\mathbb{P}(S_{t+1} = \mathbf{s}_i | S_t = \mathbf{s}_j) = p_{i,j} := \begin{cases} \pi_i & \text{if } i = j \\ (1 - \pi_i) & \text{if } i = j + 1, \\ 0 & \text{otherwise} \end{cases}$$

and initial state $S_0 = \mathbf{s}_0$. Here, S_t is the value at iteration t of the loop. Note that since $\pi_i \in \{0, 1\}$, the $p_{i,j}$ values are “degenerate” probabilities. However, this framework lets us generalize to a smooth version of the while loop naturally. To illustrate the framework, if the while loop stops at $T = 3$, the transition probabilities can be cast as a matrix

$$\mathbf{P} := (p_{i,j})_{i,j=0}^T := \begin{matrix} & \mathbf{s}_0 & \mathbf{s}_1 & \mathbf{s}_2 & \mathbf{s}_3 \\ \begin{matrix} \mathbf{s}_0 \\ \mathbf{s}_1 \\ \mathbf{s}_2 \\ \mathbf{s}_3 \end{matrix} & \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix} \end{matrix}.$$

The output \mathbf{r} of the while-loop is determined by the time at which the state stays at the same value

$$I = \min\{i \in \{1, 2, \dots\} \text{ s.t. } S_i = S_{i-1}\}.$$

Note that I itself is a random variable, as it is defined by the S_i variables. It is called a **stopping time**. The output of the chain is then

$$\begin{aligned} \mathbf{r} &= \mathbb{E}[S_I] \\ &= \sum_{i=1}^{+\infty} \mathbb{P}(I = i) \mathbb{E}[S_i | I = i] \\ &= \sum_{i=1}^{+\infty} \mathbb{P}(I = i) \mathbf{s}_{i-1} \\ &= \sum_{i=1}^{+\infty} \prod_{j=0}^{i-2} (1 - \pi_j) \pi_{i-1} \mathbf{s}_{i-1} \\ &= \sum_{i=0}^{+\infty} \prod_{j=0}^{i-1} (1 - \pi_j) \pi_i \mathbf{s}_i. \end{aligned}$$

Because the stopping time is not known ahead of time, the sum over i goes from 0 to ∞ . However, if we enforce in the stopping criterion that the while loop runs no longer than T iterations, by setting

$$\pi_i := \text{or}(f(\mathbf{s}_i), \text{eq}(i, T)) \in \{0, 1\},$$

we then naturally recover the expression found by unrolling the while loop before,

$$\mathbf{r} = \mathbb{E}[S_I] = \sum_{i=0}^T \prod_{j=0}^{i-1} (1 - \pi_j) \pi_i \mathbf{s}_i.$$

For example, with $T = 3$, the transition probability matrix is

$$\mathbf{P} = \begin{matrix} & \mathbf{s}_0 & \mathbf{s}_1 & \mathbf{s}_2 & \mathbf{s}_3 \\ \begin{matrix} \mathbf{s}_0 \\ \mathbf{s}_1 \\ \mathbf{s}_2 \\ \mathbf{s}_3 \end{matrix} & \begin{pmatrix} \pi_0 & 1 - \pi_0 & 0 & 0 \\ 0 & \pi_1 & 1 - \pi_1 & 0 \\ 0 & 0 & \pi_2 & 1 - \pi_2 \\ 0 & 0 & 0 & 1 \end{pmatrix} \end{matrix}.$$

Smoothed while loops

With the help of this framework, we can backpropagate even through the while loop's stopping criterion, provided that we smooth out the predicate. For example, we saw that the stopping criterion in Example 5.5 is $f(\mathbf{s}_i) = \text{step}(\tau - \varepsilon(\mathbf{s}_i))$ and therefore

$$\pi_i := \text{or}(f(\mathbf{s}_i), \text{eq}(i, T)) \in \{0, 1\}.$$

Due to the step function, the derivative of the while loop with respect to τ will always be 0, just like it was the case for if-else statements. If we change the stopping criterion to $f(\mathbf{s}_i) = \text{sigmoid}(\tau - \varepsilon(\mathbf{s}_i))$, we then have (recall that or is well defined on $[0, 1] \times [0, 1]$)

$$\pi_i := \text{or}(f(\mathbf{s}_i), \text{eq}(i, T)) \in [0, 1].$$

With sigmoid, we obtain more informative derivatives. In particular, with $\text{sigmoid} = \text{logistic}$, the derivatives w.r.t. τ are always non-zero. The smoothed output is expressed as before as the expectation

$$\begin{aligned} \mathbf{r} = \mathbb{E}[S_I] &= \sum_{i=0}^T \prod_{j=0}^{i-1} (1 - \pi_j) \pi_i \mathbf{s}_i \\ &= \sum_{i=0}^T \prod_{j=0}^{i-1} (1 - \text{sigmoid}(\varepsilon(\mathbf{s}_j) - \tau)) \text{sigmoid}(\varepsilon(\mathbf{s}_i) - \tau) \mathbf{s}_i. \end{aligned}$$

Instead of enforcing a number T of iterations, it is also possible to stop when the probability of stopping becomes high enough (Petersen *et al.*, 2021), assuming that the probability of stopping converges to 1.

5.11 Summary

- For conditionals, we saw that differentiating through the branch variables is not problematic given a fixed predicate.
- However, for the predicate variable, we saw that a differentiable relaxation is required to avoid null derivatives.
- We introduced soft comparison operators in a principled manner, using a stochastic process perspective, as well as the continuous extension of logical operators.

- For loops and scan define valid computational graphs, as their number of iterations is fixed ahead of time. Feedforward networks and RNNs can be seen as parameterized for loops and scan, respectively.
- Unlike for loops and scan, the number of iterations of while loops is not known ahead of time and may even be infinite. However, unrolled while loops define valid directed acyclic graphs. We defined a principled way to differentiate through the stopping criterion of a while loop, thanks to a Markov chain perspective.

6

Data structures

In computer science, a data structure is a specialized format for organizing, storing and accessing data. Mathematically, a data structure forms a so-called algebraic structure: it consists of a set and the functions to operate on that set. In this chapter, we review how to incorporate data structures into differentiable programs, with a focus on lists and dictionaries.

6.1 Lists

A list is an ordered sequence of elements. We restrict ourselves to lists whose elements all belong to the same value space \mathcal{V} . Formally, we denote a list of fixed length K with values in \mathcal{V} by a K -tuple

$$\mathbf{l} := (l_1, \dots, l_K) \in \mathcal{L}_K(\mathcal{V})$$

where each $l_i \in \mathcal{V}$ and where

$$\mathcal{L}_K(\mathcal{V}) := \mathcal{V}^K = \underbrace{\mathcal{V} \times \dots \times \mathcal{V}}_{K \text{ times}}.$$

6.1.1 Basic operations

Getting values

We first present how to retrieve values from a list $\mathbf{l} \in \mathcal{L}_K(\mathcal{V})$. We define the function $\text{list.get}: \mathcal{L}_K(\mathcal{V}) \times [K] \rightarrow \mathcal{V}$ as

$$\text{list.get}(\mathbf{l}, i) := l_i.$$

The function is continuous and differentiable in $\mathbf{l} \in \mathcal{L}_K(\mathcal{V})$ but not in $i \in [K]$, as it is a discrete variable. In the particular case $\mathcal{V} = \mathbb{R}$, $\mathcal{L}_K(\mathcal{V})$ is equivalent to \mathbb{R}^K and we can therefore write

$$\text{list.get}(\mathbf{l}, i) = \langle \mathbf{l}, \mathbf{e}_i \rangle,$$

where $\{\mathbf{e}_1, \dots, \mathbf{e}_K\}$ is the standard basis of \mathbb{R}^K .

Setting values

We now present how to replace values from a list $\mathbf{l} \in \mathcal{L}_K(\mathcal{V})$. We define the function $\text{list.set}: \mathcal{L}_K(\mathcal{V}) \times [K] \times \mathcal{V} \rightarrow \mathcal{L}_K(\mathcal{V})$ as

$$[\text{list.set}(\mathbf{l}, i, \mathbf{v})]_j := \begin{cases} \mathbf{v} & \text{if } i = j \\ l_j & \text{if } i \neq j \end{cases},$$

for $j \in [K]$. In the functional programming spirit, the function returns the **whole** new list, even though a single element has been modified. Again, the function is continuous and differentiable in $\mathbf{l} \in \mathcal{L}_K(\mathcal{V})$ and $\mathbf{v} \in \mathcal{V}$ but not in $i \in [K]$. In the particular case $\mathcal{V} = \mathbb{R}$, given a list $\mathbf{l} = (l_1, \dots, l_K)$, we can write

$$\text{list.set}(\mathbf{l}, i, v) = (v - l_i)\mathbf{e}_i.$$

That is, we subtract the old value l_i and add the new value v at the location $i \in [K]$.

Implementation

A fixed-length list can be implemented as an **array**, which enables $O(1)$ random access to individual elements. The hardware counterpart of an array is random access memory (RAM), in which memory can be retrieved by address (location).

6.1.2 Operations on variable-length lists

So far, we focused on lists of fixed length K . We now turn our attention to variable-length lists, whose size can decrease or increase over time. In addition to the `list.get` and `list.set` functions, they support functions that can change the size of a list.

Initializing lists

In order to initialize a list, we define `list.init`: $\mathcal{V} \rightarrow \mathcal{L}_1(\mathcal{V})$ as

$$\text{list.init}(\mathbf{v}) := (\mathbf{v}),$$

where we used (\mathbf{v}) to denote a 1-tuple.

Pushing values

In order to add new values either to the left or to the right, we define `list.pushLeft`: $\mathcal{L}_K(\mathcal{V}) \times \mathcal{V} \rightarrow \mathcal{L}_{K+1}(\mathcal{V})$ as

$$\text{list.pushLeft}(\mathbf{l}, \mathbf{v}) := (\mathbf{v}, \mathbf{l}_1, \dots, \mathbf{l}_K).$$

and `list.pushRight`: $\mathcal{L}_K(\mathcal{V}) \times \mathcal{V} \rightarrow \mathcal{L}_{K+1}(\mathcal{V})$ as

$$\text{list.pushRight}(\mathbf{l}, \mathbf{v}) := (\mathbf{l}_1, \dots, \mathbf{l}_K, \mathbf{v}).$$

Popping values

In order to remove values either from the left or from the right, we define `list.popLeft`: $\mathcal{L}_K(\mathcal{V}) \rightarrow \mathcal{L}_{K-1}(\mathcal{V}) \times \mathcal{V}$ as

$$\text{list.popLeft}(\mathbf{l}) := (\mathbf{l}_2, \dots, \mathbf{l}_K), \mathbf{l}_1$$

and `list.popRight`: $\mathcal{L}_K(\mathcal{V}) \rightarrow \mathcal{L}_{K-1}(\mathcal{V}) \times \mathcal{V}$ as

$$\text{list.popRight}(\mathbf{l}) := (\mathbf{l}_1, \dots, \mathbf{l}_{K-1}), \mathbf{l}_K.$$

The set $\mathcal{L}_0(\mathcal{V})$ is a singleton which contains the empty list.

Inserting values

The `pushLeft` and `pushRight` functions can only insert values at the beginning and at the end of a list, respectively. We now study the `insert` function, whose goal is to be able to add a new value at an arbitrary location, shifting all values to the right and increasing the list size by 1. We define the function $\text{list.insert} : \mathcal{L}_K(\mathcal{V}) \times [K+1] \times \mathcal{V} \rightarrow \mathcal{L}_{K+1}(\mathcal{V})$ as

$$[\text{list.insert}(\mathbf{l}, i, \mathbf{v})]_j := \begin{cases} l_j & \text{if } j < i \\ \mathbf{v} & \text{if } j = i, \\ l_{j-1} & \text{if } j > i \end{cases}$$

for $j \in [K+1]$. As for the `list.set` function, `list.insert` is readily continuous and differentiable in \mathbf{l} and \mathbf{v} , but not in i , as it is a discrete variable. As special cases, we naturally recover

$$\begin{aligned} \text{list.insert}(\mathbf{l}, 1, \mathbf{v}) &= \text{pushLeft}(\mathbf{l}, \mathbf{v}), \\ \text{list.insert}(\mathbf{l}, K+1, \mathbf{v}) &= \text{pushRight}(\mathbf{l}, \mathbf{v}). \end{aligned}$$

Differentiability

The `list.init`, `list.push` and `list.pop` functions are readily continuous and differentiable with respect to their arguments (a continuous relaxation is not needed). As for the `list.set` function, the `list.insert` function is continuous and differentiable in \mathbf{l} and \mathbf{v} , but not in i .

Implementation

Under the hood, a variable-length list can be implemented as a linked list or as a dynamic array. A linked list gives $O(K)$ random access while a dynamic array allows $O(1)$ random access, at the cost of memory reallocations.

Stacks and queues

The `list.pushRight` and `list.popRight` functions can be used to implement a **stack** (last in first out a.k.a. LIFO behavior). The `list.pushLeft` and `list.popRight` functions can be used to implement a **queue** (first in first out a.k.a. FIFO behavior).

6.1.3 Continuous relaxations using soft indexing

Getting values

In order to be able to differentiate `list.get` w.r.t. indexing, a natural idea is to replace the integer index $i \in [K]$ by a distribution $\pi_i \in \Delta^K$, which we can interpret as a **soft index**. An integer index $i \in [K]$ is then equivalent to a **delta distribution** $\pi_i \in \{e_1, \dots, e_K\}$. We define the continuous relaxation `list.softGet`: $\mathcal{L}_K(\mathcal{V}) \times \Delta^K \rightarrow \text{conv}(\mathcal{V})$ as

$$\begin{aligned} \text{list.softGet}(\mathbf{l}, \pi_i) &:= \sum_{j=1}^K \pi_{i,j} \mathbf{l}_j \\ &= \text{cond}(\pi_i, \mathbf{l}_1, \dots, \mathbf{l}_K) \\ &= \mathbb{E}_{I \sim \text{Categorical}(\pi_i)}[\mathbf{l}_I], \end{aligned}$$

where `cond` is studied in Section 5.7. In the particular case $\mathcal{V} = \mathbb{R}$, we obtain

$$\text{list.softGet}(\mathbf{l}, i) = \langle \mathbf{l}, \pi_i \rangle.$$

This is illustrated in Fig. 6.1.

The choice of the distribution $\pi_i = (\pi_{i,1}, \dots, \pi_{i,K})$ encodes the importance of the elements $(\mathbf{l}_1, \dots, \mathbf{l}_K)$ w.r.t. \mathbf{l}_i . If we consider that the smaller $|i - j|$ is, the more related \mathbf{l}_i and \mathbf{l}_j are, then it makes sense to define a distribution centered around i (i.e., such that the mode of the distribution is achieved at i). For example, limiting ourselves to the **neighbors** \mathbf{l}_{i-1} and \mathbf{l}_{i+1} (i.e., a window of size 1), we can define the sparse distribution

$$\pi_i := \frac{1}{4} \cdot e_{i-1} + \frac{1}{2} e_i + \frac{1}{4} \cdot e_{i+1} \in \Delta^K.$$

In this particular case, the continuous relaxation of the `list.get` function can then be expressed as a **discrete convolution**,

$$\text{list.softGet}(\mathbf{l}, \pi_i) = (\text{list.get}(\mathbf{l}, \cdot) * \kappa)(i) = \sum_{j=-\infty}^{\infty} \text{list.get}(\mathbf{l}, i - j) \kappa(j),$$

where $\kappa(-1) := \frac{1}{4}$, $\kappa(1) := \frac{1}{4}$, $\kappa(0) := \frac{1}{2}$, and $\kappa(j) := 0$ for $j \notin \{-1, 0, 1\}$. Assuming $\mathcal{V} = \mathbb{R}^M$, the computational complexity of `list.softGet` is $O(M \cdot |\text{supp}(\pi_i)|)$.

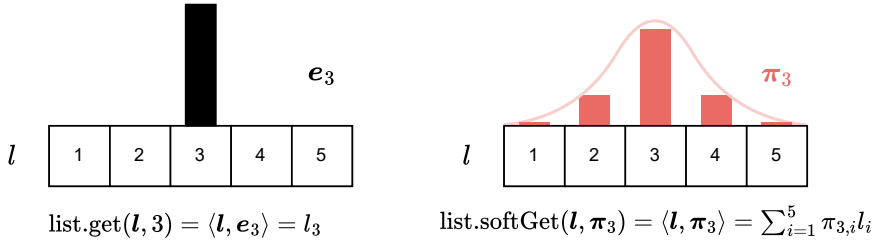


Figure 6.1: The $\text{list.get}(l, i)$ function is continuous and differentiable in l but not in i . Its relaxation $\text{list.softGet}(l, \pi_i)$ is differentiable in both l and π_i . When $\mathcal{V} = \mathbb{R}$, $\text{list.softGet}(l, \pi_i)$ can be seen as taking the inner product between the list l and the probability distribution π_i , instead of the delta distribution (canonical vector) e_i .

Setting values

To differentiate w.r.t. indexing, we can define the continuous relaxation $\text{list.softSet}: \mathcal{L}_K(\mathcal{V}) \times \triangle^K \times \mathcal{V} \rightarrow \mathcal{L}_K(\text{conv}(\mathcal{V}))$ as

$$\begin{aligned} [\text{list.softSet}(l, \pi_i, v)]_j &:= \mathbb{E}[\text{list.set}(l, I, v)]_j \\ &= \mathbb{P}(I = j)v + \mathbb{P}(I \neq j)l_j \\ &= \pi_{i,j}v + (1 - \pi_{i,j})l_j, \end{aligned}$$

where $j \in [K]$ and $I \sim \text{Categorical}(\pi_i)$. Equivalently, we can write

$$\begin{aligned} \text{list.softSet}(l, \pi_i, v) &= (\pi_{i,1}v + (1 - \pi_{i,1})l_1, \dots, \pi_{i,K}v + (1 - \pi_{i,K})l_K) \\ &= (\text{ifelse}(\pi_{i,1}, v, l_1), \dots, \text{ifelse}(\pi_{i,K}, v, l_K)), \end{aligned}$$

where ifelse is studied in Section 5.6. Since

$$\text{ifelse}(\pi, u_1, u_0) = \mathbb{E}_{I \sim \text{Bernoulli}(\pi)}[u_I],$$

this relaxation amounts to using an element-wise expectation. As a result, the list output by list.softSet takes values in $\text{conv}(\mathcal{V})$ instead of \mathcal{V} . Note however that when $\mathcal{V} = \mathbb{R}^M$, then $\text{conv}(\mathcal{V}) = \mathbb{R}^M$ as well.

Inserting values

To differentiate value insertion w.r.t. indexing, we can define the continuous relaxation $\text{list.softInsert} : \mathcal{L}_K(\mathcal{V}) \times \triangle^{K+1} \times \mathcal{V} \rightarrow \mathcal{L}_{K+1}(\text{conv}(\mathcal{V}))$

$$\begin{aligned} [\text{list.softInsert}(\mathbf{l}, \boldsymbol{\pi}_i, \mathbf{v})]_j &:= \mathbb{E}[\text{list.insert}(\mathbf{l}, I, \mathbf{v})] \\ &= \mathbb{P}(I > j)\mathbf{l}_j + \mathbb{P}(I = j)\mathbf{v} + \mathbb{P}(I < j)\mathbf{l}_{j-1}, \end{aligned}$$

where $I \sim \text{Categorical}(\boldsymbol{\pi}_i)$. The three necessary probabilities can easily be calculated for $j \in [K + 1]$ by

$$\begin{aligned} \mathbb{P}(I > j) &= \begin{cases} 0 & \text{if } j = K + 1 \\ \sum_{k=j+1}^{K+1} \pi_{i,k} & \text{otherwise} \end{cases} \\ \mathbb{P}(I = j) &= \pi_{i,j} \\ \mathbb{P}(I < j) &= \begin{cases} 0 & \text{if } j = 1 \\ \sum_{k=1}^{j-1} \pi_{i,k} & \text{otherwise} \end{cases}. \end{aligned}$$

Multi-dimensional indexing

In multi-dimensional lists (arrays or tensors), each element $\mathbf{l}_i \in \mathcal{V}$ of a list $\mathbf{l} \in \mathcal{L}_{K_1, \dots, K_T}(\mathcal{V})$ can now be indexed by a multivariate integer $\mathbf{i} = (i_1, \dots, i_T) \in [K_1] \times \dots \times [K_T]$, where $T \in \mathbb{N}$ is the number of axes of \mathbf{l} . We can always flatten a multi-dimensional list into an uni-dimensional list by replacing the multi-dimensional index $\mathbf{i} \in [K_1] \times \dots \times [K_T]$ by a flat index $i \in [K_1 \dots K_T]$. The converse operation, converting a flat uni-dimensional array into a multi-dimensional array, is also possible. Therefore, there is a **bijection** between $[K]$ and $[K_1] \times \dots \times [K_T]$ for $K := K_1 \dots K_T$.

This means that the previous discussion on soft indexing in the uni-dimensional setting readily applies to the multi-dimensional setting. All it takes is the ability to define a probability distribution $\boldsymbol{\pi}_i \in \triangle^{K_1 \times \dots \times K_T}$. For example, when working with images, we can define a probability distribution putting probability mass only on the neighboring pixels of pixel \mathbf{i} , a standard approach in image processing. Another simple approach is to use a product of axis-wise probability distributions.

6.2 Dictionaries

A dictionary (a.k.a. associative array or map) is an unordered list of **key-value pairs**, such that each possible key appears at most once in the list. We denote the set of keys by \mathcal{K} and the set of values by \mathcal{V} (both being potentially infinite). We can then define the set of dictionaries of size L from \mathcal{K} to \mathcal{V} by

$$\mathcal{D}_L(\mathcal{K}, \mathcal{V}) := \mathcal{L}_L(\mathcal{K} \times \mathcal{V}) = (\mathcal{K} \times \mathcal{V})^L$$

and one such dictionary by

$$\mathbf{d} := ((\mathbf{k}_1, \mathbf{v}_1), \dots, (\mathbf{k}_L, \mathbf{v}_L)) \in \mathcal{D}_L(\mathcal{K}, \mathcal{V}).$$

6.2.1 Basic operations

Getting values

The goal of the `dict.get` function is to retrieve the value associated with a key, assuming that the dictionary contains this key. Formally, we define the `dict.get`: $\mathcal{D}_L(\mathcal{K}, \mathcal{V}) \times \mathcal{K} \rightarrow \mathcal{V} \cup \{\infty\}$ function as

$$\text{dict.get}(\mathbf{d}, \mathbf{k}) := \begin{cases} \mathbf{v}_i & \text{if } \exists i \in [L] \text{ s.t. } \mathbf{k} = \mathbf{k}_i \\ \infty & \text{if } \mathbf{k} \notin \{\mathbf{k}_1, \dots, \mathbf{k}_L\} \end{cases}.$$

The function is continuous and differentiable in the dictionary \mathbf{d} , but not in the key \mathbf{k} . Equivalently, we can write the function as

$$\text{dict.get}(\mathbf{d}, \mathbf{k}) := \frac{\sum_{i=1}^L \text{eq}(\mathbf{k}, \mathbf{k}_i) \mathbf{v}_i}{\sum_{i=1}^L \text{eq}(\mathbf{k}, \mathbf{k}_i)}.$$

The denominator encodes the fact that the function is undefined if no key in the dictionary \mathbf{d} matches the key \mathbf{k} . Assuming $\mathbf{k} \in \{\mathbf{k}_1, \dots, \mathbf{k}_L\}$ and $\mathcal{V} = \mathbb{R}^M$, we can also write

$$\text{dict.get}(\mathbf{d}, \mathbf{k}) = \mathbf{v}_i \quad \text{where} \quad i = \arg \max_{j \in [L]} \|\mathbf{k} - \mathbf{k}_j\|_2,$$

which shows that we can see `dict.get` as a **nearest neighbor search**.

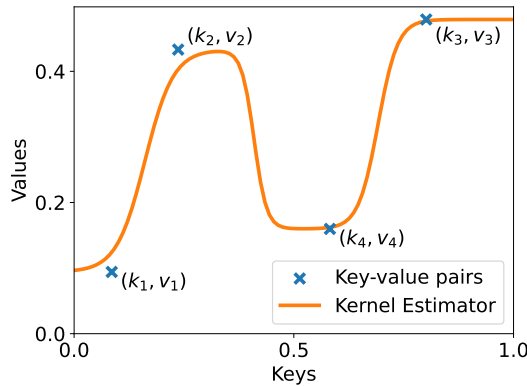


Figure 6.2: Given a set of key-value pairs $(\mathbf{k}_i, \mathbf{v}_i) \in \mathcal{K} \times \mathcal{V}$ defining a dictionary \mathbf{d} , we can estimate a continuous mapping from \mathcal{K} to \mathcal{V} using Nadaraya–Watson kernel regression (here, illustrated with $\mathcal{K} = \mathcal{V} = \mathbb{R}$). When keys are normalized to have unit norm, this recovers softargmax attention from Transformers.

Setting values

The goal of the `dict.set` function is to replace the value associated with an existing key. Formally, we define the $\text{dict.set} : \mathcal{D}_L(\mathcal{K}, \mathcal{V}) \times \mathcal{K} \times \mathcal{V} \rightarrow \mathcal{D}_L(\mathcal{K}, \mathcal{V})$ function as

$$(\text{dict.set}(\mathbf{d}, \mathbf{k}, \mathbf{v}))_i := \begin{cases} (\mathbf{k}_i, \mathbf{v}) & \text{if } \mathbf{k}_i = \mathbf{k} \\ (\mathbf{k}_i, \mathbf{v}_i) & \text{if } \mathbf{k}_i \neq \mathbf{k} \end{cases}.$$

The function leaves the dictionary unchanged if no key in the dictionary matches the input key \mathbf{k} . The function is continuous and differentiable in \mathbf{d} and \mathbf{v} , but not in \mathbf{k} .

Implementation

While we view dictionaries as lists of key-value pairs, in practice, a dictionary (a.k.a. associative array) is often implemented using a hash table or search trees. The hardware counterpart of a dictionary is called content-addressable memory (CAM), a.k.a. associative memory.

6.2.2 Continuous relaxation using kernel regression

A dictionary can be seen as a (potentially non-injective) function that associates a value \mathbf{v} to each key \mathbf{k} . To obtain a continuous relaxation of the operations associated to a dictionary, we can adopt a probabilistic perspective of the mapping from keys to values. We can view keys and values as two continuous random variables K and V . We can express the conditional PDF $f(\mathbf{v}|\mathbf{k})$ of $V|K$ in terms of the joint PDF $f(\mathbf{k}, \mathbf{v})$ of (K, V) and the marginal PDF $f(\mathbf{k})$ of K as

$$f(\mathbf{v}|\mathbf{k}) = \frac{f(\mathbf{k}, \mathbf{v})}{f(\mathbf{k})}.$$

Integrating, we obtain the **conditional expectation**

$$\mathbb{E}[V|K = \mathbf{k}] = \int_{\mathcal{V}} f(\mathbf{v}|\mathbf{k}) \mathbf{v} d\mathbf{v} = \int_{\mathcal{V}} \frac{f(\mathbf{k}, \mathbf{v})}{f(\mathbf{k})} \mathbf{v} d\mathbf{v}.$$

This is the **Bayes predictor**, in the sense that $\mathbb{E}[V|K]$ is the minimizer of $\mathbb{E}[(h(K) - V)^2]$ over the space of measurable functions $h: \mathcal{K} \rightarrow \mathcal{V}$. Using a sample of L input-output pairs $(\mathbf{k}_i, \mathbf{v}_i)$, corresponding to key-value pairs in our case, **Nadaraya–Watson kernel regression** estimates the joint PDF and the marginal PDF using **kernel density estimation** (KDE). Using a product of isotropic kernels κ_σ and ρ_σ for key-value pairs, we can define

$$\hat{f}_\sigma(\mathbf{k}, \mathbf{v}) := \frac{1}{L} \sum_{i=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_i) \rho_\sigma(\mathbf{v} - \mathbf{v}_i).$$

The corresponding marginal distribution on the keys is then given as

$$\begin{aligned} \hat{f}_\sigma(\mathbf{k}) &:= \int_{\mathcal{V}} \hat{f}_\sigma(\mathbf{k}, \mathbf{v}) d\mathbf{v} \\ &= \frac{1}{L} \sum_{i=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_i) \int_{\mathcal{V}} \rho_\sigma(\mathbf{v} - \mathbf{v}_i) d\mathbf{v} \\ &= \frac{1}{L} \sum_{i=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_i). \end{aligned}$$

Replacing f with \hat{f}_σ , we obtain the following estimator of the conditional expectation

$$\begin{aligned}\hat{\mathbb{E}}[V|K = \mathbf{k}] &:= \int_{\mathcal{V}} \frac{\hat{f}_\sigma(\mathbf{k}, \mathbf{v})}{\hat{f}_\sigma(\mathbf{k})} \mathbf{v} d\mathbf{v} \\ &= \int_{\mathcal{V}} \frac{\frac{1}{L} \sum_{i=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_i) \rho_\sigma(\mathbf{v} - \mathbf{v}_i)}{\frac{1}{L} \sum_{i=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_i)} \mathbf{v} d\mathbf{v} \\ &= \frac{\sum_{i=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_i) \int_{\mathcal{V}} \rho_\sigma(\mathbf{v} - \mathbf{v}_i) \mathbf{v} d\mathbf{v}}{\sum_{i=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_i)} \\ &= \frac{\sum_{i=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_i) \mathbf{v}_i}{\sum_{i=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_i)}.\end{aligned}$$

In the above, we assumed that $\rho_\sigma(\mathbf{v} - \mathbf{v}_i) = p_{\mathbf{v}_i, \sigma}(\mathbf{v})$, where $p_{\mathbf{v}_i, \sigma}(\mathbf{v})$ is the PDF of a distribution whose mean is \mathbf{v}_i , so that

$$\int_{\mathcal{V}} \rho_\sigma(\mathbf{v} - \mathbf{v}_i) \mathbf{v} d\mathbf{v} = \mathbb{E}_{V \sim p_{\mathbf{v}_i, \sigma}}[V] = \mathbf{v}_i.$$

Given a dictionary $\mathbf{d} = ((\mathbf{k}_1, \mathbf{v}_1), \dots, (\mathbf{k}_L, \mathbf{v}_L))$, we can therefore define the `dict.softGet`: $\mathcal{D}_L(\mathcal{K}, \mathcal{V}) \times \mathcal{K} \rightarrow \text{conv}(\mathcal{V})$ function as

$$\text{dict.softGet}(\mathbf{d}, \mathbf{k}) := \frac{\sum_{i=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_i) \mathbf{v}_i}{\sum_{i=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_i)}.$$

This kernel regression perspective on dictionaries was previously pointed out by Zhang *et al.* (2021). It is illustrated in Fig. 6.2 with $\mathcal{K} = \mathcal{V} = \mathbb{R}$.

6.2.3 Discrete probability distribution perspective

While the set of possible keys \mathcal{K} is potentially infinite, the set of keys $\{\mathbf{k}_1, \dots, \mathbf{k}_L\} \subset \mathcal{K}$ associated with a particular dictionary $\mathbf{d} = ((\mathbf{k}_1, \mathbf{v}_1), \dots, (\mathbf{k}_L, \mathbf{v}_L))$ is finite. To a particular key \mathbf{k} , we can therefore associate a discrete probability distribution $\boldsymbol{\pi}_{\mathbf{k}} = (\pi_{\mathbf{k},1}, \dots, \pi_{\mathbf{k},L}) \in \Delta^L$ over the keys $(\mathbf{k}_1, \dots, \mathbf{k}_L)$ of \mathbf{d} , defined by

$$\pi_{\mathbf{k},i} := \frac{\kappa_\sigma(\mathbf{k} - \mathbf{k}_i)}{\sum_{j=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_j)} \quad \forall i \in [L].$$

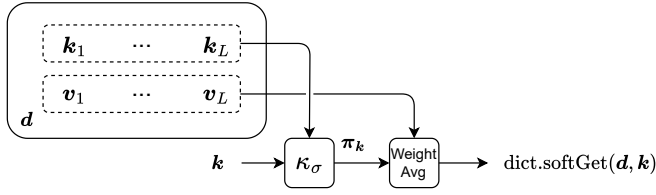


Figure 6.3: Computation graph of the `dict.softGet` function. We can use a kernel κ_σ to produce a discrete probability distribution $\pi_k = (\pi_{k,1}, \dots, \pi_{k,L}) \in \Delta^L$, that captures the affinity between the dictionary keys (k_1, \dots, k_L) and the input key k . The `dict.softGet` function can then merely be seen as a convex combination (weighted average) of values (v_1, \dots, v_L) using the probability values $(\pi_{k,1}, \dots, \pi_{k,L})$ as weights.

This distribution captures the affinity between the input key k and the keys (k_1, \dots, k_L) of dictionary d . As illustrated in Fig. 6.3, we obtain

$$\begin{aligned} \text{dict.softGet}(d, k) &= \mathbb{E}_{i \sim \text{Categorical}(\pi_k)}[v_i] \\ &= \sum_{i=1}^L \pi_{k,i} v_i. \end{aligned}$$

In the limit $\sigma \rightarrow 0$, we recover

$$\text{dict.get}(d, k) = \frac{\sum_{i=1}^L \text{eq}(k, k_i) v_i}{\sum_{i=1}^L \text{eq}(k, k_i)}.$$

While the `dict.get` function is using a mapping from keys $k \in \{k_1, \dots, k_L\}$ to integer indices $[L]$, the `dict.softGet` function is using a mapping from keys $k \in \{k_1, \dots, k_L\}$ to distributions $\pi_k \in \Delta^L$. This perspective allows us to reuse the soft functions we developed for lists in Section 6.1. For example, we can softly replace the value associated with key k by performing

$$\text{list.softSet}(d, \pi_k, (k, v)).$$

Unlike `dict.set`, the function is differentiable w.r.t. the distribution π_k .

6.2.4 Link with attention in Transformers

In the case when κ_σ is the Gaussian kernel, assuming that the keys are normalized to have unit norm (which is often the case in practical

implementations (Schlag *et al.*, 2021; Dehghani *et al.*, 2023)), we obtain

$$\begin{aligned}\kappa_\sigma(\mathbf{k} - \mathbf{k}_i) &= \exp(-\|\mathbf{k} - \mathbf{k}_i\|_2^2 / (2\sigma^2)) \\ &= \exp(-(\|\mathbf{k}\|_2^2 + \|\mathbf{k}_i\|_2^2) / (2\sigma^2)) \exp(\langle \mathbf{k}, \mathbf{k}_i \rangle / \sigma^2) \\ &= \exp(-\sigma^2) \exp(\langle \mathbf{k}, \mathbf{k}_i \rangle / \sigma^2)\end{aligned}$$

so that

$$\begin{aligned}\pi_{\mathbf{k},i} &= \frac{\kappa_\sigma(\mathbf{k} - \mathbf{k}_i)}{\sum_{j=1}^L \kappa_\sigma(\mathbf{k} - \mathbf{k}_j)} \\ &= \frac{\exp(\langle \mathbf{k}, \mathbf{k}_i \rangle / \sigma^2)}{\sum_{j=1}^L \exp(\langle \mathbf{k}, \mathbf{k}_j \rangle / \sigma^2)}.\end{aligned}$$

We recognize the softargmax operator. Given, a dictionary $\mathbf{d} = ((\mathbf{k}_1, \mathbf{v}_1), \dots, (\mathbf{k}_L, \mathbf{v}_L))$, we thus recover attention from Transformers (Vaswani *et al.*, 2017) as

$$\text{dict.softGet}(\mathbf{d}, \mathbf{k}) = \frac{\exp(\langle \mathbf{k}, \mathbf{k}_i \rangle / \sigma^2) \mathbf{v}_i}{\sum_{j=1}^L \exp(\langle \mathbf{k}, \mathbf{k}_j \rangle / \sigma^2)}.$$

Transformers can therefore be interpreted as relying on a differentiable dictionary mechanism. Besides Transformers, content-based memory addressing is also used in neural Turing machines (Graves *et al.*, 2014).

6.3 Summary

- Operations on lists are continuous and differentiable w.r.t. the list, but not w.r.t. the integer index. Similarly, operations on dictionaries are continuous and differentiable w.r.t. the dictionary, but not w.r.t. the input key.
- Similarly to the way we handled the predicate in conditionals, we can replace the integer index (respectively the key) with a probability distribution over the indices (respectively the keys).
- This allows us to obtain a probabilistic relaxation of operations on lists. In particular, the relaxation for list.get amounts to performing a convolution. The relaxation for dict.get amounts to computing a conditional expectation using kernel regression.

- When using a Gaussian kernel with keys normalized to unit norm, we recover softargmax attention from Transformers.

Part III

Differentiating through programs

7

Finite differences

One of the simplest way to numerically compute derivatives is to use finite differences, which approximate the infinitesimal definition of derivatives. Finite differences only require **function evaluations**, and can therefore work with blackbox functions (i.e., they ignore the compositional structure of functions). Without loss of generality, our exposition focuses on computing directional derivatives $\partial f(\mathbf{w})[\mathbf{v}]$, for a function $f: \mathcal{E} \rightarrow \mathcal{F}$, evaluated at $\mathbf{w} \in \mathcal{E}$, in the direction $\mathbf{v} \in \mathcal{E}$.

7.1 Forward differences

From Definition 2.4 and Definition 2.13, the directional derivative and more generally the JVP are defined as a limit,

$$\partial f(\mathbf{w})[\mathbf{v}] := \lim_{\delta \rightarrow 0} \frac{f(\mathbf{w} + \delta \mathbf{v}) - f(\mathbf{w})}{\delta}.$$

This suggests that we can approximate the directional derivative and the JVP using

$$\partial f(\mathbf{w})[\mathbf{v}] \approx \frac{f(\mathbf{w} + \delta \mathbf{v}) - f(\mathbf{w})}{\delta},$$

for some $0 < \delta \ll 1$. This formula is called a **forward difference**. From the Taylor expansion in Section 2.5.4, we indeed have

$$f(\mathbf{w} + \delta \mathbf{v}) - f(\mathbf{w}) = \delta \partial f(\mathbf{w})[\mathbf{v}] + \frac{\delta^2}{2} \partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{v}] + \frac{\delta^3}{3!} \partial^3 f(\mathbf{w})[\mathbf{v}, \mathbf{v}, \mathbf{v}] + \dots$$

so that

$$\begin{aligned} \frac{f(\mathbf{w} + \delta \mathbf{v}) - f(\mathbf{w})}{\delta} &= \partial f(\mathbf{w})[\mathbf{v}] + \frac{\delta}{2} \partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{v}] + \frac{\delta^2}{3!} \partial^3 f(\mathbf{w})[\mathbf{v}, \mathbf{v}, \mathbf{v}] + \dots \\ &= \partial f(\mathbf{w})[\mathbf{v}] + o(\delta). \end{aligned}$$

The error incurred by choosing a finite rather than infinitesimal δ in the forward difference formula is called the **truncation error**. The Taylor approximation above shows that this error is of the order of $o(\delta)$.

However, we cannot choose a too small value of δ , because the evaluation of the function f on a computer rounds the value of f to machine precision. Mathematically, a scalar-valued function f evaluated on a computer becomes a function \tilde{f} such that $\tilde{f}(w) \approx [f(w)/\varepsilon]\varepsilon$, where $[f(w)/\varepsilon]$ denotes the closest integer of $f(w)/\varepsilon \in \mathbb{R}$ and ε is the machine precision, i.e., the smallest non-zero real number encoded by the machine. This means that the difference $f(\mathbf{w} + \delta \mathbf{v}) - f(\mathbf{w})$ evaluated on a computer is prone to **round-off error** of the order of $o(\varepsilon)$. We illustrate the trade-off between truncation and round-off errors in Fig. 7.1.

7.2 Backward differences

As an alternative, we can approximate the directional derivative and the JVP by

$$\partial f(\mathbf{w})[\mathbf{v}] \approx \frac{f(\mathbf{w}) - f(\mathbf{w} - \delta \mathbf{v})}{\delta},$$

for some $0 < \delta \ll 1$. This formula is called a **backward difference**. From the Taylor expansion in Section 2.5.4, we easily verify that $(f(\mathbf{w}) - f(\mathbf{w} - \delta \mathbf{v}))/\delta = \partial f(\mathbf{w})[\mathbf{v}] + o(\delta)$, so that the truncation error is the same as for the forward difference.

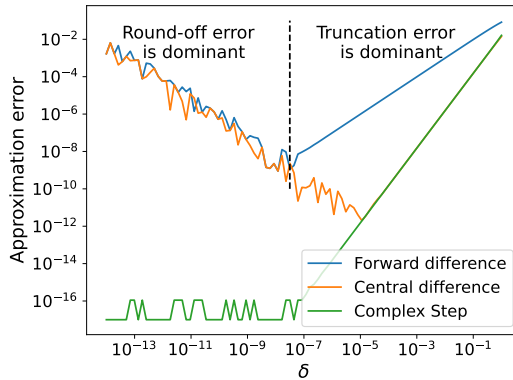


Figure 7.1: Numerical differentiation of $f(x) := \text{softplus}(x) = \log(1 + \exp(x))$, to approximate $f'(x) = \text{logistic}(x)$ at $x = 1$. The forward and central difference methods induce both truncation error (for large δ) and round-off error (for small δ). The complex step method enjoys smaller round-off error.

7.3 Central differences

Rather than using an asymmetric formula to approximate the derivative, as in forward and backward differences, we can use a symmetric formula

$$\partial f(\mathbf{w})[\mathbf{v}] \approx \frac{f(\mathbf{w} + \delta \mathbf{v}) - f(\mathbf{w} - \delta \mathbf{v})}{2\delta},$$

for some $0 < \delta \ll 1$. This formula is called a **central difference**. From the Taylor expansion in Section 2.5.4, we have

$$\begin{aligned} f(\mathbf{w} + \delta \mathbf{v}) - f(\mathbf{w} - \delta \mathbf{v}) &= 2\delta \partial f(\mathbf{w})[\mathbf{v}] + \frac{2\delta^3}{3!} \partial^3 f(\mathbf{w})[\mathbf{v}, \mathbf{v}, \mathbf{v}] \\ &\quad + \frac{2\delta^5}{5!} \partial^5 f(\mathbf{w})[\mathbf{v}, \dots, \mathbf{v}] + \dots \end{aligned}$$

so that

$$\begin{aligned} \frac{f(\mathbf{w} + \delta \mathbf{v}) - f(\mathbf{w} - \delta \mathbf{v})}{2\delta} &= \partial f(\mathbf{w})[\mathbf{v}] + \frac{\delta^2}{3!} \partial^3 f(\mathbf{w})[\mathbf{v}, \mathbf{v}, \mathbf{v}] + \dots \\ &= \partial f(\mathbf{w})[\mathbf{v}] + o(\delta^2). \end{aligned}$$

We see that the terms corresponding to derivatives of **even order** canceled out, allowing the formula to achieve $o(\delta^2)$ truncation error.

For any $\delta < 1$, the truncation error of the central difference is much smaller than the one of the forward or backward differences as confirmed empirically in Fig. 7.1.

7.4 Higher-accuracy finite differences

The truncation error can be further reduced by making use of additional function evaluations. One can generalize the forward difference scheme by a formula of the form

$$\partial f(\mathbf{w})[\mathbf{v}] \approx \sum_{i=0}^p \frac{a_i}{\delta} f(\mathbf{w} + i\delta\mathbf{v})$$

requiring $p+1$ evaluations. To select the a_i and reach a truncation error of order $o(\delta^p)$, we can use a Taylor expansion on each term of the sum to get

$$\sum_{k=0}^p \frac{a_i}{\delta} f(\mathbf{w} + i\delta\mathbf{v}) = \sum_{k=0}^p a_i \sum_{j=0}^p \frac{i^j \delta^{j-1}}{j!} \partial^j f(\mathbf{w})[\mathbf{v}, \dots, \mathbf{v}] + o(\delta^p).$$

By grouping the terms in the sum for each order of derivative, we obtain a set of $p+1$ equations to be satisfied by the $p+1$ coefficients a_0, \dots, a_p , that is,

$$\begin{aligned} a_0 + a_1 + \dots + a_p &= 0 \\ a_1 + 2a_2 + \dots + pa_p &= 1 \\ a_1 + 2^j a_2 + \dots + p^j a_p &= 0 \quad \forall j \in \{2, \dots, p\}. \end{aligned}$$

This system of equations can be solved analytically to derive the coefficients. Backward differences can be generalized similarly by using $\partial f(\mathbf{w})[\mathbf{v}] \approx \sum_{i=0}^p \frac{a_i}{\delta} f(\mathbf{w} - i\delta\mathbf{v})$. Similarly, the central difference scheme can be generalized by using

$$\partial f(\mathbf{w})[\mathbf{v}] \approx \sum_{i=-p}^p \frac{a_i}{\delta} f(\mathbf{w} + i\delta\mathbf{v}),$$

to reach a truncation error of order $o(\delta^{2p})$. Solving for the coefficients a_{-p}, \dots, a_p as above reveals that $a_0 = 0$. Therefore, only $2p$ evaluations are necessary.

7.5 Higher-order finite differences

To approximate higher order derivatives, we can follow a similar reasoning. Namely, we can generalize the forward difference scheme to approximate the derivative of order k by

$$\partial^k f(\mathbf{w})[\mathbf{v}, \dots, \mathbf{v}] \approx \sum_{i=0}^p \frac{a_i}{\delta^k} f(\mathbf{w} + i\delta\mathbf{v}).$$

As before, we can expand the terms in the sum. For the approximation to capture only the k^{th} derivative, we now require the coefficients a_i to satisfy

$$\begin{aligned} 0^j a_0 + 1^j a_1 + 2^j a_2 + \dots + p^j a_p &= 0 \quad \forall j \in \{0, \dots, k-1\}. \\ 0^k a_0 + 1^k a_1 + 2^k a_2 + \dots + p^k a_p &= k! \\ 0^j a_0 + 1^j a_1 + 2^j a_2 + \dots + p^j a_p &= 0 \quad \forall j \in \{k+1, \dots, p\}. \end{aligned}$$

With the resulting coefficients, we obtain a truncation error of order $o(\delta^{p-k+1})$, while making $p+1$ evaluations. For example, for $p = k = 2$, we can approximate the second-order derivative as

$$\partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{v}] \approx \frac{-(3/2)f(\mathbf{x}) + 2f(\mathbf{x} + \delta\mathbf{v}) - (1/2)f(\mathbf{x} + 2\delta\mathbf{v})}{\delta^2},$$

with a truncation error of order $o(\delta)$.

The central difference scheme can be generalized similarly by

$$\partial^k f(\mathbf{w})[\mathbf{v}, \dots, \mathbf{v}] \approx \sum_{i=-p}^p \frac{a_i}{\delta^k} f(\mathbf{w} + i\delta\mathbf{v}),$$

to reach truncation errors of order $o(\delta^{2p+2-2\lceil(k+1)/2\rceil})$. For example, for $k = 2$, $p = 1$, we obtain the second-order central difference

$$\partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{v}] \approx \frac{f(\mathbf{w} + \delta\mathbf{v}) + f(\mathbf{w} - \delta\mathbf{v}) - 2f(\mathbf{w})}{\delta^2}.$$

By using a Taylor expansion we see that, this time, the terms corresponding to derivatives of **odd order** cancel out and the truncation error is $o(\delta^2)$ while requiring 3 evaluations.

7.6 Complex-step derivatives

Suppose f is well defined on \mathbb{C}^P , the space of P -dimensional complex numbers. Let us denote the imaginary unit by $i = \sqrt{-1}$. Then, the Taylor expansion of f reads

$$\begin{aligned} f(\mathbf{w} + (i\delta)\mathbf{v}) &= f(\mathbf{w}) + (i\delta)\partial f(\mathbf{w})[\mathbf{v}] + \frac{(i\delta)^2}{2}\partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{v}] \\ &\quad + \frac{(i\delta)^3}{3!}\partial^3 f(\mathbf{w})[\mathbf{v}, \mathbf{v}, \mathbf{v}] + \dots \\ &= f(\mathbf{w}) + (i\delta)\partial f(\mathbf{w})[\mathbf{v}] - \frac{\delta^2}{2}\partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{v}] \\ &\quad - \frac{i(\delta)^3}{3!}\partial^3 f(\mathbf{w})[\mathbf{v}, \mathbf{v}, \mathbf{v}] + \dots \end{aligned}$$

We see that the real part corresponds to even-degree terms and the imaginary part corresponds to odd-degree terms. We therefore obtain

$$\operatorname{Re}(f(\mathbf{w} + (i\delta)\mathbf{v})) = f(\mathbf{w}) + o(\delta^2)$$

and

$$\operatorname{Im}\left(\frac{f(\mathbf{w} + (i\delta)\mathbf{v})}{\delta}\right) = \partial f(\mathbf{w})[\mathbf{v}] + o(\delta^2).$$

This suggests that we can compute directional derivatives using the approximation

$$\partial f(\mathbf{w})[\mathbf{v}] \approx \operatorname{Im}\left(\frac{f(\mathbf{w} + (i\delta)\mathbf{v})}{\delta}\right),$$

for $0 < \delta \ll 1$. This is called the **complex-step derivative** approximation (Squire and Trapp, 1998; Martins *et al.*, 2003).

Contrary to forward, backward and central differences, we see that only a **single function call** is necessary. A function call on complex numbers may take roughly twice the cost of a function call on real numbers. However, thanks to the fact that a difference of functions is no longer needed, the complex-step derivative approximation usually enjoys smaller round-off error as illustrated in Fig. 7.1. That said, one drawback of the method is that all elementary operations within the program implementing the function f must be well-defined on complex numbers, e.g., using overloading.

Table 7.1: Computational complexity in number of function evaluations for computing the directional derivative and the gradient of a function $f: \mathbb{R}^P \rightarrow \mathbb{R}$ by finite differences and complex step derivatives.

	Directional derivative	Gradient
Forward difference	2	$P + 1$
Backward difference	2	$P + 1$
Central difference	2	$2P$
Complex step	1	P

7.7 Complexity

We now discuss the computational complexity in terms of function evaluations of finite differences and complex-step derivatives. For concreteness, as this is the most common use case in machine learning, we discuss the case of a single $M = 1$ output, i.e., we want to differentiate a function $f: \mathbb{R}^P \rightarrow \mathbb{R}$. Whether we use forward, backward or central differences, the computational complexity of computing the directional derivative $\partial f(\mathbf{w})[\mathbf{v}]$ in any direction \mathbf{v} amounts to two calls to f . For computing the gradient $\nabla f(\mathbf{w})$, we can use (see Definition 2.7) that

$$[\nabla f(\mathbf{w})]_j = \langle \nabla f(\mathbf{w}), \mathbf{e}_j \rangle = \partial f(\mathbf{w})[\mathbf{e}_j],$$

for $j \in [P]$. For forward and backward differences, we therefore need $P + 1$ function calls to compute the gradient, while we need $2P$ function calls for central differences. For the complex step approximation, we need P complex function calls. We summarize the complexities in Table 7.1.

7.8 Summary

- Finite differences are a simple way to numerically compute derivatives using only function evaluations.
- Central differences achieve smaller truncation error than forward and backward differences. It is possible to achieve smaller truncation error, at the cost of more function evaluations.

- Complex-step derivatives achieve smaller round-off error than central differences but require the function and the program implementing it to be well-defined on complex numbers.
- However, whatever the method used, finite differences require a number of function calls that is proportional to the number of dimensions. They are therefore seldom used in machine learning, where there can be millions or billions of dimensions. The main use cases of finite differences are therefore i) for blackbox functions of low dimension and ii) for test purposes (e.g., checking that a gradient function is correctly implemented).
- For modern machine learning, the main workhorse is automatic differentiation, as it leverages the compositional structure of functions. This is what we study in the next chapter.

8

Automatic differentiation

In Chapter 2, we reviewed the fundamentals of differentiation and stressed the importance of two linear maps: the Jacobian-vector product (JVP) and its adjoint, the vector-Jacobian product (VJP). In this chapter, we review **forward-mode** and **reverse-mode** autodiff using these two linear maps. We start with **computation chains** and then generalize to feedforward networks and general **computation graphs**. We also review checkpointing, reversible layers and randomized estimators.

8.1 Computation chains

To begin with, consider a **computation chain** (Section 4.1.1) representing a function $f: \mathcal{S}_0 \rightarrow \mathcal{S}_K$ expressed as a sequence of compositions $f := f_K \circ \dots \circ f_1$, where $f_k: \mathcal{S}_{k-1} \rightarrow \mathcal{S}_k$. The computation of f can be

unrolled into a sequence of operations

$$\begin{aligned}
 \mathbf{s}_0 &\in \mathcal{S}_0 \\
 \mathbf{s}_1 &:= f_1(\mathbf{s}_0) \in \mathcal{S}_1 \\
 &\vdots \\
 \mathbf{s}_K &:= f_K(\mathbf{s}_{K-1}) \in \mathcal{S}_K \\
 f(\mathbf{s}_0) &:= \mathbf{s}_K.
 \end{aligned} \tag{8.1}$$

Our goal is to compute the variations of f around a given input \mathbf{s}_0 . In a feedforward network, this amounts to estimating the influence of a given input \mathbf{s}_0 for fixed parameters (we will see how to estimate the variations w.r.t. parameters \mathbf{w} in the sequel).

Jacobian matrix. We first consider the computation of the full Jacobian $\partial f(\mathbf{s}_0)$, seen as a **matrix**, as the notation ∂ indicates. Following Proposition 2.2, we have

$$\partial f(\mathbf{s}_0) = \partial f_K(\mathbf{s}_{K-1}) \dots \partial f_1(\mathbf{s}_0), \tag{8.2}$$

where $\partial f_k(\mathbf{s}_{k-1})$ are the Jacobians of the intermediate functions computed at $\mathbf{s}_0, \dots, \mathbf{s}_K$, as defined in Eq. (8.1). We may also want to compute the transpose of the Jacobian

$$\partial f(\mathbf{s}_0)^\top = \partial f_1(\mathbf{s}_0)^\top \dots \partial f_K(\mathbf{s}_{K-1})^\top.$$

In both cases, the main drawback of this approach is computational: computing the full $\partial f(\mathbf{s}_0)$ requires to materialize the intermediate Jacobians in memory and to perform matrix-matrix multiplications. However, in practice, computing the full Jacobian is rarely needed. Indeed, oftentimes, we only need to right-multiply or left-multiply with $\partial f(\mathbf{s}_0)$. This gives rise to forward-mode and reverse-mode autodiff, respectively.

8.1.1 Forward-mode

We now interpret the Jacobian $\partial f(\mathbf{s}_0)$ as a **linear map**, as the non-bold ∂ indicates. Following Proposition 2.6, $\partial f(\mathbf{s}_0)$ is the composition of the

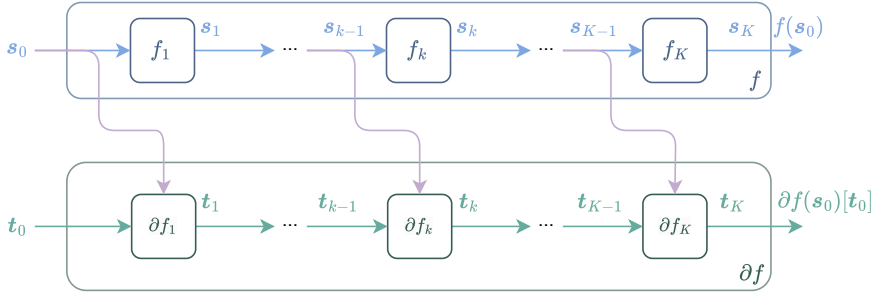


Figure 8.1: Forward-mode autodiff for a computation chain. For readability, we denoted the intermediate JVP as a function of two variables $\partial f_k : s_{k-1}, t_{k-1} \mapsto \partial f_k(s_{k-1})[t_{k-1}]$ with $\partial f_k(s_{k-1})[t_{k-1}] = t_k$.

intermediate linear maps,

$$\partial f(s_0) = \partial f_K(s_{K-1}) \circ \dots \circ \partial f_1(s_0).$$

Evaluating $\partial f(s_0)$ on an **input direction** $v \in \mathcal{S}_0$ can be decomposed, like the function Eq. (8.1) itself, into intermediate computations

$$\begin{aligned} t_0 &:= v \\ t_1 &:= \partial f_1(s_0)[t_0] \\ &\vdots \\ t_K &:= \partial f_K(s_{K-1})[t_{K-1}] \\ \partial f(s_0)[v] &:= t_K. \end{aligned}$$

Each intermediate $\partial f_k(s_{k-1})[t_{k-1}]$ amounts to a Jacobian-vector product (JVP) and can be performed in a **forward** manner, along the computation of the intermediate states s_k . This can also be seen as multiplying the matrix defined in Eq. (8.2) with a vector, from **right to left**. This is illustrated in Fig. 8.1 and the procedure is summarized in Algorithm 8.1.

Computational complexity. The JVP follows exactly the computations of f , with an additional variable t_k being propagated. If we consider that computing ∂f_k is roughly as costly as computing f_k , then computing a

Algorithm 8.1 Forward-mode autodiff for computation chains

Functions: $f := f_K \circ \dots \circ f_1$ **Inputs:** input $s_0 \in \mathcal{S}_0$, input direction $v \in \mathcal{S}_0$ 1: Initialize $t_0 := v$ 2: **for** $k := 1, \dots, K$ **do**3: Compute $s_k := f_k(s_{k-1}) \in \mathcal{S}_k$ 4: Compute $t_k := \partial f_k(s_{k-1})[t_{k-1}] \in \mathcal{S}_k$ **Outputs:** $f(s_0) := s_K, \partial f(s_0)[v] = t_K$

JVP has roughly twice the computational cost of f . See Section 8.3.3 for a more general and more formal statement.

Memory usage. The memory usage of a program at a given evaluation step is the number of variables that need to be stored in memory to ensure the execution of all remaining steps. The memory cost of a program is then the maximal memory usage over all evaluation steps. For our purposes, we analyze the memory usage and memory cost by examining the given program. Formal definitions of operations on memory such as read, write, delete and associated memory costs are presented by Griewank and Walther (2008, Chapter 4).

For example, to execute the chain $f = f_K \circ \dots \circ f_1$, at each step k , we only need to have access to s_{k-1} to execute the rest of the program. As we compute s_k , we can delete s_{k-1} from memory and replace it by s_k . Therefore, the memory cost associated to the evaluation of f is just the maximal dimension of the s_k variables.

For forward mode autodiff, as we follow the computations of f , at each step k , we only need to have access to s_{k-1} and t_{k-1} to execute the rest of the program. The memory used by s_{k-1} and t_{k-1} can directly be used for s_k, t_k once they are computed. The memory usage associated to the JVP is summarized in Fig. 8.2. Overall the memory cost of the JVP is then exactly twice the memory cost of the function itself.

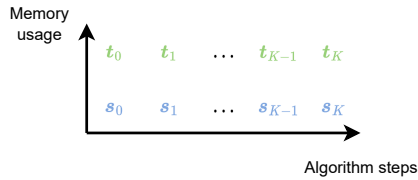


Figure 8.2: Memory usage of forward-mode autodiff for a computation chain. Here $t_0 = v$, $s_K = f(s_0)$, $t_K = \partial f(s_0)[v]$.

8.1.2 Reverse-mode

In machine learning, most functions whose gradient we need to compute take the form $\ell \circ f$, where ℓ is a scalar-valued loss function and f is a network. As seen in Proposition 2.3, the gradient takes the form

$$\nabla(\ell \circ f)(s_0) = \partial f(s_0)^*[\nabla\ell(f(s_0))].$$

This motivates the need for applying the **adjoint** $\partial f(s_0)^*$ to $\nabla\ell(f(s_0)) \in \mathcal{S}_K$ and more generally to any **output direction** $u \in \mathcal{S}_K$. From Proposition 2.7, we have

$$\partial f(s_0)^* = \partial f_1(s_0)^* \circ \dots \circ \partial f_K(s_{K-1})^*.$$

Evaluating $\partial f(s_0)^*$ on an output direction $u \in \mathcal{S}_K$ is decomposed as

$$\begin{aligned} r_K &:= u \\ r_{K-1} &:= \partial f_K(s_{K-1})^*[r_K] \\ &\vdots \\ r_0 &:= \partial f_1(s_0)^*[r_1] \\ \partial f(s_0)^*[u] &:= r_0. \end{aligned}$$

Each intermediate adjoint $\partial f_k(s_{k-1})^*$ amounts to a vector-Jacobian product (VJP). The key difference with the forward mode is that the procedure runs **backward** through the chain, hence the name **reverse mode** autodiff. This can also be seen as multiplying Eq. (8.2) from **left to right**. The procedure is illustrated in Fig. 8.3 and summarized in Algorithm 8.2.

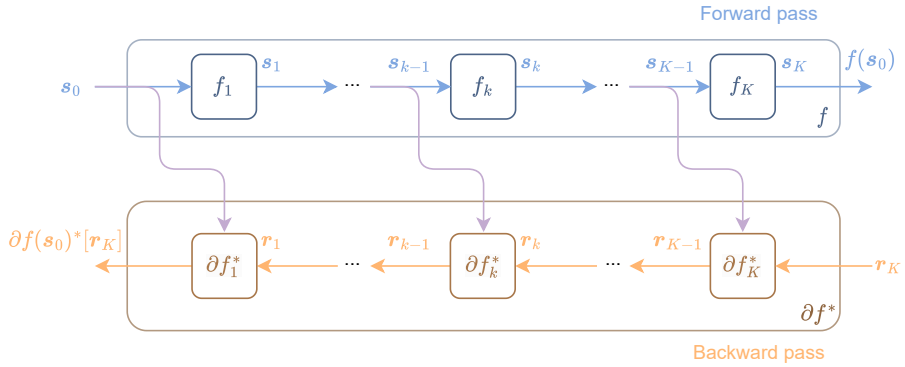


Figure 8.3: Reverse mode of automatic differentiation for a computation chain. For readability, we denoted the intermediate VJPs as functions of two variables $\partial f_k^* : (s_{k-1}, r_k) \mapsto \partial f_k(s_{k-1})^*[r_k]$, with $\partial f_k(s_{k-1})^*[r_k] = r_{k-1}$.

Algorithm 8.2 Reverse-mode autodiff for computation chains

Functions: $f := f_K \circ \dots \circ f_1$,

Inputs: input $s_0 \in \mathcal{S}_0$, output direction $u \in \mathcal{S}_K$

- 1: **for** $k := 1, \dots, K$ **do** ▷ Forward pass
- 2: Compute $s_k := f_k(s_{k-1}) \in \mathcal{S}_k$
- 3: Initialize $r_K := u$.
- 4: **for** $k := K, \dots, 1$ **do** ▷ Backward pass
- 5: Compute $r_{k-1} := \partial f_k(s_{k-1})^*[r_k] \in \mathcal{S}_{k-1}$

Outputs: $f(s_0) := s_K$, $\partial f(s_0)^*[u] = r_0$

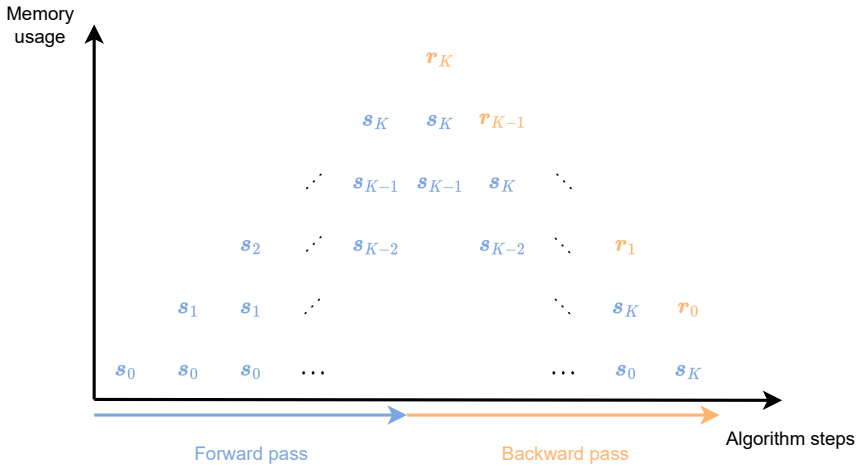


Figure 8.4: Memory usage of reverse mode autodiff for a computation chain.

Computational complexity. In terms of number of operations, the VJP simply passes two times through the chain, once forward, then backward. If we consider the intermediate VJPs to be roughly as costly as the intermediate functions themselves, the VJP amounts just to twice the cost of the original function, just as the JVP. See Section 8.3.3 for a more generic and formal statement.

Memory usage. Recall that the memory usage of a program at a given evaluation step is the number of variables that need to be stored in memory to ensure the execution of the remaining steps. If we inspect Algorithm 8.3, to execute all backward steps, that is the loop in line 4, we need to have access to all the intermediate inputs s_0, \dots, s_{K-1} . Therefore, the memory cost of reverse-mode autodiff is proportional to the length of the chain K . Fig. 8.4 illustrates the memory usage during reverse mode autodiff. It grows linearly until the end of the forward pass and then progressively decreases until it outputs the value of the function and the VJP. The memory cost can be mitigated by means of checkpointing techniques presented in Section 8.5.

Decoupled function and VJP evaluations. The additional memory cost of reverse mode autodiff comes with some advantages. If we need to compute $\partial f(\mathbf{s}_0)^*[\mathbf{u}_i]$ for n different output directions \mathbf{u}_i , we only need to compute and store once the intermediate computations \mathbf{s}_k and then make n calls to the backward pass. In other words, by storing in memory the intermediate computations \mathbf{s}_k , we may instantiate a **VJP operator**, which we may apply to any \mathbf{u} through the backward pass. Formally, the forward and backward passes can be decoupled as

$$\text{forward}(f, \mathbf{s}_0) := (f(\mathbf{s}_0), \partial f(\mathbf{s}_0)^*)$$

where

$$\partial f(\mathbf{s}_0)^*[\mathbf{u}] := \text{backward}(\mathbf{u}; \mathbf{s}_0, \dots, \mathbf{s}_{K-1}).$$

In functional programming terminology, the VJP $\partial f(\mathbf{s}_0)^*$ is a **closure**, as it contains the intermediate computations $\mathbf{s}_0, \dots, \mathbf{s}_K$. The same can be done for the JVP $\partial f(\mathbf{s}_0)$ if we want to apply to multiple directions \mathbf{v}_i .

Example 8.1 (Multilayer perceptron with fixed parameters). As a running example, consider a multilayer perceptron (MLP) with one hidden layer and (for now) given fixed weights. As presented in Chapter 4, an MLP can be decomposed as

$$\begin{aligned} \mathbf{s}_0 &= \mathbf{x} \\ \mathbf{s}_1 &= f_1(\mathbf{s}_0) = \sigma(\mathbf{A}_1 \mathbf{s}_0 + \mathbf{b}_1) \\ \mathbf{s}_2 &= f_2(\mathbf{s}_1) = \mathbf{A}_2 \mathbf{s}_1 + \mathbf{b}_2 \\ f(\mathbf{x}) &= \mathbf{s}_2, \end{aligned}$$

for $\mathbf{A}_1, \mathbf{A}_2, \mathbf{b}_1, \mathbf{b}_2$ some fixed parameters and σ an activation function such as the softplus activation function $\sigma(x) = \log(1 + e^x)$ with derivative $\sigma'(x) = e^x / (1 + e^x)$.

Evaluating the JVP of f on an input \mathbf{x} along a direction \mathbf{v} can

then be decomposed as

$$\begin{aligned} \mathbf{t}_0 &= \mathbf{v} \\ \mathbf{t}_1 &= \sigma'(\mathbf{A}_1 \mathbf{s}_0 + \mathbf{b}_1) \odot (\mathbf{A}_1 \mathbf{t}_0) \\ \mathbf{t}_2 &= \mathbf{A}_2 \mathbf{t}_1 \\ \partial f(\mathbf{x})[\mathbf{v}] &= \mathbf{t}_2, \end{aligned}$$

where we used in the second line the JVP of element-wise function as in Example 8.4.

Evaluating the VJP of f at \mathbf{x} requires to evaluate the intermediate VJPs at the stored activations

$$\begin{aligned} \mathbf{r}_2 &= \mathbf{u} \\ \mathbf{r}_1 &= \partial f_2(\mathbf{s}_1)^*[\mathbf{r}_2] = \mathbf{A}_2^\top \mathbf{r}_2 \\ \mathbf{r}_0 &= \partial f_1(\mathbf{s}_0)^*[\mathbf{r}_1] = \mathbf{A}_1^\top (\sigma'(\mathbf{A}_1 \mathbf{s}_0 + \mathbf{b}_1) \odot \mathbf{r}_1) \\ \partial f(\mathbf{x})^*[\mathbf{u}] &= \mathbf{r}_0. \end{aligned}$$

8.1.3 Complexity of computing entire Jacobians

In this section, we analyze the time and space complexities of forward-mode and reverse-mode autodiff for computing the **entire** Jacobian matrix $\partial f(\mathbf{s}_0)$ of a computation chain $f = f_K \circ \dots \circ f_1$, where $f_k: \mathcal{S}_{k-1} \rightarrow \mathcal{S}_k$. We assume $\mathcal{S}_k \subseteq \mathbb{R}^{D_k}$, $D_K = M$ and $D_0 = D$. Therefore, we have $f: \mathbb{R}^D \rightarrow \mathbb{R}^M$ and $\partial f(\mathbf{s}_0) \in \mathbb{R}^{M \times D}$.

Complexity of forward-mode autodiff

Using Definition 2.9, we find that we can extract each **column** $[\partial f(\mathbf{s}_0)]_{:,j} \in \mathbb{R}^M$ of the Jacobian matrix, for $j \in [D]$, by multiplying with the standard basis vector $\mathbf{e}_j \in \mathbb{R}^D$:

$$\begin{aligned} [\partial f(\mathbf{s}_0)]_{:,1} &= \partial f(\mathbf{s}_0)[\mathbf{e}_1] \\ &\vdots \\ [\partial f(\mathbf{s}_0)]_{:,D} &= \partial f(\mathbf{s}_0)[\mathbf{e}_D]. \end{aligned}$$

Computing the full Jacobian matrix therefore requires D JVPs with vectors in \mathbb{R}^D . Assuming each f_k in the chain composition has the form

$f_k : \mathbb{R}^{D_{k-1}} \rightarrow \mathbb{R}^{D_k}$, seen as a matrix, $\partial f_k(\mathbf{s}_{k-1})$ has size $D_k \times D_{k-1}$. Therefore, the computational cost of D JVPs is $O\left(D \sum_{k=1}^K D_k D_{k-1}\right)$. The memory cost is $O(\max_{k \in [K]} D_k)$, since we can release intermediate computations after each layer is processed. Setting $D_1 = \dots = D_{K-1} = D$ for simplicity and using $D_K = M$, we obtain that the computational cost of computing D JVPs and therefore of computing the full Jacobian matrix by forward-mode autodiff is $O(MD^2 + KD^3)$. The memory cost is $O(\max\{D, M\})$. If a function has a single-input ($D = 1$), then the forward mode computes the entire Jacobian at once, which reduces to a single directional derivative.

Complexity of reverse-mode autodiff

Using Definition 2.9, we find that we can extract each **row** of the Jacobian matrix $[\partial f(\mathbf{s}_0)]_i \in \mathbb{R}^D$, for $i \in [M]$, by multiplying with the standard basis vector $\mathbf{e}_i \in \mathbb{R}^M$:

$$\begin{aligned} [\partial f(\mathbf{s}_0)]_1 &= \partial f(\mathbf{s}_0)^*[\mathbf{e}_1] \\ &\vdots \\ [\partial f(\mathbf{s}_0)]_M &= \partial f(\mathbf{s}_0)^*[\mathbf{e}_M]. \end{aligned}$$

Computing the full Jacobian matrix therefore requires M VJPs with vectors in \mathbb{R}^M . Assuming as before that each f_k in the chain composition has the form $f_k : \mathbb{R}^{D_{k-1}} \rightarrow \mathbb{R}^{D_k}$, the computational cost of M VJPs is $O\left(M \sum_{k=1}^K D_k D_{k-1}\right)$. However, the memory cost is $O(\sum_{k=1}^K D_k)$, as we need to store the intermediate computations for each of the K layers. Setting $D_0 = \dots = D_{K-1} = D$ for simplicity and using $D_K = M$, we obtain that the computational cost of computing M VJPs and therefore of computing the full Jacobian matrix by reverse-mode autodiff is $O(M^2D + KMD^2)$. The memory cost is $O(KD + M)$. If the function has a single output ($M = 1$), reverse-mode autodiff computes the entire Jacobian at once, which reduces to the gradient.

When to use forward-mode vs. reverse-mode autodiff?

We summarize the time and space complexities in Table 8.1. Generally, if $M < D$, reverse-mode is more advantageous at the price of some

	Forward-mode	Reverse-mode
Time	$O(MD^2 + KD^3)$	$O(M^2D + KMD^2)$
Space	$O(\max\{M, D\})$	$O(KD + M)$

Table 8.1: Time and space complexities of forward-mode and reverse-mode autodiff for computing the full Jacobian of a chain of functions $f = f_K \circ \dots \circ f_1$, where $f_k: \mathbb{R}^D \rightarrow \mathbb{R}^D$ if $k = 1, \dots, K-1$ and $f_K: \mathbb{R}^D \rightarrow \mathbb{R}^M$. We assume ∂f_k is a dense linear operator. Forward mode requires D JVPs. Reverse mode requires M VJP.

memory cost. If $M \geq D$, forward mode is more advantageous.

8.2 Feedforward networks

In the previous section, we derived forward-mode autodiff and reverse-mode autodiff for computation chains with an input $\mathbf{s}_0 \in \mathcal{S}_0$. In this section, we now derive reverse-mode autodiff for feedforward networks, in which each layer f_k is now allowed to depend explicitly on some additional parameters $\mathbf{w}_k \in \mathcal{W}_k$. The recursion is

$$\begin{aligned}
 \mathbf{s}_0 &:= \mathbf{x} \in \mathcal{S}_0 \\
 \mathbf{s}_1 &:= f_1(\mathbf{s}_0, \mathbf{w}_1) \in \mathcal{S}_1 \\
 &\vdots \\
 \mathbf{s}_K &:= f_K(\mathbf{s}_{K-1}, \mathbf{w}_K) \in \mathcal{S}_K \\
 f(\mathbf{x}, \mathbf{w}) &:= \mathbf{s}_K,
 \end{aligned}$$

where $\mathcal{S}_0 = \mathcal{X}$ and $\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_K) \in \mathcal{W}_1 \times \dots \times \mathcal{W}_K$. Each f_k is now a function of two arguments. The first argument depends on the previous layer, but the second argument does not. This is illustrated in Fig. 8.5. We now explain how to differentiate a feedforward network.

8.2.1 Computing the adjoint

The function has the form $f: \mathcal{E} \rightarrow \mathcal{F}$, where $\mathcal{E} := \mathcal{X} \times (\mathcal{W}_1 \times \dots \times \mathcal{W}_K)$ and $\mathcal{F} := \mathcal{S}_K$. From Section 2.3, we know that the VJP has the form $\partial f(\mathbf{x}, \mathbf{w})^*: \mathcal{F} \rightarrow \mathcal{E}$. Therefore, we want to be able to compute $\partial f(\mathbf{x}, \mathbf{w})^*[\mathbf{u}] \in \mathcal{E}$ for any $\mathbf{u} \in \mathcal{F}$.

Fortunately, the backward recursion is only a slight modification of the computation chain case. Indeed, since $f_k: \mathcal{E}_k \rightarrow \mathcal{F}_k$, where $\mathcal{E}_k := \mathcal{S}_{k-1} \times \mathcal{W}_k$ and $\mathcal{F}_k := \mathcal{S}_k$, the intermediate VJPs have the form $\partial f_k(\mathbf{s}_{k-1}, \mathbf{w}_k)^*: \mathcal{F}_k \rightarrow \mathcal{E}_k$. We therefore arrive at the recursion

$$\begin{aligned} \mathbf{r}_K &= \mathbf{u} \in \mathcal{S}_K \\ (\mathbf{r}_{K-1}, \mathbf{g}_K) &= \partial f_K(\mathbf{s}_{K-1}, \mathbf{w}_K)^*[\mathbf{r}_K] \in \mathcal{S}_{K-1} \times \mathcal{W}_K \\ &\vdots \\ (\mathbf{r}_0, \mathbf{g}_1) &= \partial f_1(\mathbf{s}_0, \mathbf{w}_1)^*[\mathbf{r}_1] \in \mathcal{S}_0 \times \mathcal{W}_1. \end{aligned}$$

The final output is

$$\partial f(\mathbf{x}, \mathbf{w})^*[\mathbf{u}] = (\mathbf{r}_0, (\mathbf{g}_1, \dots, \mathbf{g}_K)).$$

8.2.2 Computing the gradient

We often compose a network with a loss function

$$L(\mathbf{w}; \mathbf{x}, \mathbf{y}) := \ell(f(\mathbf{x}, \mathbf{w}); \mathbf{y}) \in \mathbb{R}.$$

From Proposition 2.7, the gradient is given by

$$\nabla L(\mathbf{w}; \mathbf{x}, \mathbf{y}) = (\mathbf{g}_1, \dots, \mathbf{g}_K) \in \mathcal{W}_1 \times \dots \times \mathcal{W}_K$$

where

$$\partial f(\mathbf{x}, \mathbf{w})^*[\mathbf{u}] = (\mathbf{r}_0, (\mathbf{g}_1, \dots, \mathbf{g}_K)),$$

with $\mathbf{u} = \nabla \ell(f(\mathbf{x}, \mathbf{w}); \mathbf{y}) \in \mathcal{S}_K$. The output $\mathbf{r}_0 \in \mathcal{S}_0$, where $\mathcal{S}_0 = \mathcal{X}$, corresponds to the gradient w.r.t. $\mathbf{x} \in \mathcal{X}$ and is typically not needed, except in generative modeling settings. The full procedure is summarized in Algorithm 8.3.

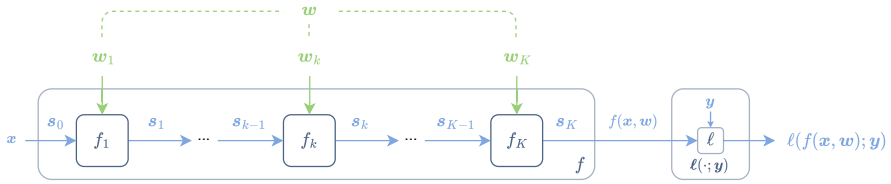


Figure 8.5: Computation graph of an MLP as a function of its parameters.

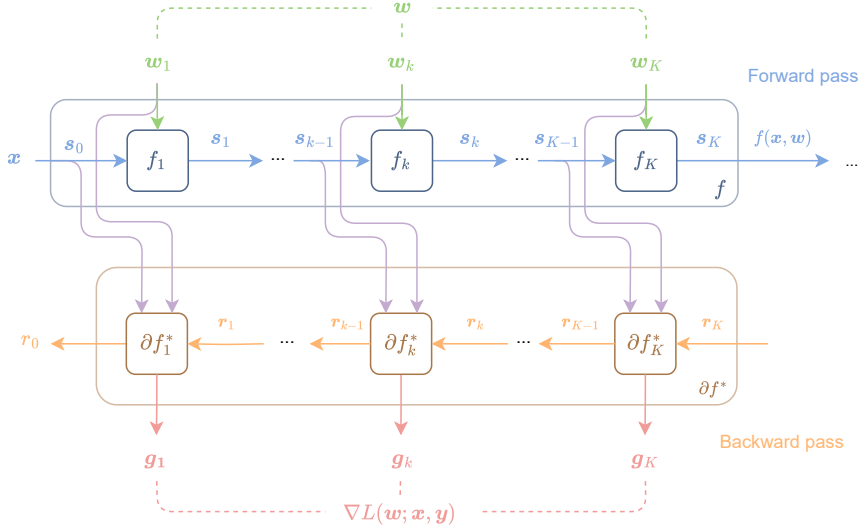


Figure 8.6: Reverse mode of automatic differentiation, a.k.a., gradient back-propagation to compute the gradient of the loss of an MLP on an input label pair. For readability, we denoted the intermediate VJPs as functions of three variables $\partial f_k^* : (s_{k-1}, w_{k-1}, r_k) \mapsto \partial f_k(s_{k-1}, w_k)^*[r_k]$ with $\partial f_k(s_{k-1}, w_k)^*[r_k] = (r_{k-1}, g_k)$.

Algorithm 8.3 Gradient back-propagation for feedforward networks

Functions: f_1, \dots, f_K in sequential order**Inputs:** data point $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$ parameters $\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_K) \in \mathcal{W}_1 \times \dots \times \mathcal{W}_K$

- 1: Initialize $\mathbf{s}_0 := \mathbf{x}$ ▷ Forward pass
 - 2: **for** $k := 1, \dots, K$ **do**
 - 3: Compute and store $\mathbf{s}_k := f_k(\mathbf{s}_{k-1}, \mathbf{w}_k) \in \mathcal{S}_k$
 - 4: Compute $\ell(\mathbf{s}_K; \mathbf{y})$ and $\mathbf{u} := \nabla \ell(\mathbf{s}_K; \mathbf{y}) \in \mathcal{S}_K$
 - 5: Initialize $\mathbf{r}_K := \mathbf{u} \in \mathcal{S}_K$ ▷ Backward pass
 - 6: **for** $k := K, \dots, 1$ **do**
 - 7: Compute $(\mathbf{r}_{k-1}, \mathbf{g}_k) := \partial f_k(\mathbf{s}_{k-1}, \mathbf{w}_k)^*[\mathbf{r}_k] \in \mathcal{S}_{k-1} \times \mathcal{W}_k$
 - 8: **Outputs:** $L(\mathbf{w}; \mathbf{x}, \mathbf{y}) := \ell(\mathbf{s}_K; \mathbf{y})$, $\nabla L(\mathbf{w}; \mathbf{x}, \mathbf{y}) = (g_1, \dots, g_K)$
-

8.3 Computation graphs

In the previous sections, we reviewed autodiff for computation chains and its extension to feedforward networks. In this section, we review its generalization to computation graphs, introduced in Section 4.1.3. Our formalism assumes, without loss of generality, that functions can take multiple inputs but only produce a single output, as we can always group multiple outputs as a tuple. This enables a one-to-one correspondence between output variables \mathbf{s}_k and functions f_k .

8.3.1 Forward mode

The forward mode corresponds to computing the JVP of the program in an input direction $\mathbf{v} \in \mathcal{S}_0$. In the case of a **computation chain** $f := f_K \circ \dots \circ f_1$, each f_k takes a **single** input $\mathbf{s}_{k-1} \in \mathcal{S}_{k-1}$ and each $\partial f_k(\mathbf{s}_{k-1})$ takes a single input direction $\mathbf{t}_{k-1} \in \mathcal{S}_{k-1}$. The forward pass on iteration $k \in [K]$ is then, starting from $\mathbf{t}_0 := \mathbf{v}$,

$$\begin{aligned}\mathbf{s}_k &:= f_k(\mathbf{s}_{k-1}) \in \mathcal{S}_k \\ \mathbf{t}_k &:= \partial f_k(\mathbf{s}_{k-1})[\mathbf{t}_{k-1}] \in \mathcal{S}_k.\end{aligned}$$

In the case of a **computation graph**, specified by functions f_1, \dots, f_K in topological order, each f_k may now take **multiple** inputs $(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}}) \in$

$\mathcal{S}_{i_1} \times \cdots \times \mathcal{S}_{i_{p_k}}$, and each $\partial f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}})$ takes as many input directions $(\mathbf{t}_{i_1}, \dots, \mathbf{t}_{i_{p_k}}) \in \mathcal{S}_{i_1} \times \cdots \times \mathcal{S}_{i_{p_k}}$, where $(i_1, \dots, i_{p_k}) := \text{pa}(k)$ are the parents of f_k . The forward pass on step $k \in [K]$ then computes both intermediate inputs and directions as

$$\begin{aligned}\mathbf{s}_k &:= f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}}) \in \mathcal{S}_k \\ \mathbf{t}_k &:= \partial f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}})[\mathbf{t}_{i_1}, \dots, \mathbf{t}_{i_{p_k}}] \in \mathcal{S}_k.\end{aligned}$$

Let us recall that $\partial_i f_k$ means that we differentiate f_k w.r.t. to its i^{th} argument. Using the fan-in rule in Proposition 2.8, we obtain that the derivatives are propagated as

$$\mathbf{t}_k = \sum_{j=1}^{p_k} \partial_j f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}})[\mathbf{t}_{i_j}] \in \mathcal{S}_k.$$

The final output is $\partial f(\mathbf{s}_0)[\mathbf{v}] = \mathbf{t}_K$. The resulting generic forward-mode procedure is summarized in Fig. 8.7 and in Algorithm 8.4. Intuitively, the algorithm consists in computing and summing intermediate JVPs along the forward pass. Although not explicitly mentioned, we can release \mathbf{s}_k and \mathbf{t}_k from memory when no child node depends on node k .

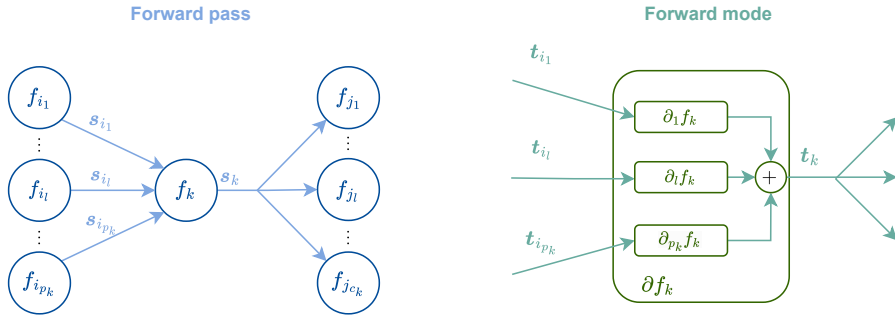


Figure 8.7: Forward mode automatic differentiation in a computation graph.

Algorithm 8.4 Forward-mode autodiff for computation graphs

Functions: f_1, \dots, f_K in topological order

Inputs: input $s_0 \in \mathcal{S}_0$, input direction $v \in \mathcal{S}_0$

- 1: Initialize $t_0 := v$
- 2: **for** $k := 1, \dots, K$ **do** ▷ Forward pass
- 3: Retrieve parent nodes $(i_1, \dots, i_{p_k}) := \text{pa}(k)$
- 4: Compute $s_k := f_k(s_{i_1}, \dots, s_{i_{p_k}}) \in \mathcal{S}_k$
- 5: Compute

$$\begin{aligned}
 t_k &:= \partial f_k(s_{i_1}, \dots, s_{i_{p_k}})[t_{i_1}, \dots, t_{i_{p_k}}] \\
 &= \sum_{j=1}^{p_k} \partial_j f_k(s_{i_1}, \dots, s_{i_{p_k}})[t_{i_j}] \in \mathcal{S}_k.
 \end{aligned}$$

- 6: **Outputs:** $f(s_0) := s_K \in \mathcal{S}_K$, $\partial f(s_0)[v] = t_K \in \mathcal{S}_K$
-

8.3.2 Reverse mode

The reverse mode corresponds to computing the VJP of the program in an output direction $\mathbf{u} \in \mathcal{S}_K$. In the case of a **computation chain** $f := f_K \circ \dots \circ f_1$, each intermediate variable \mathbf{s}_k is **used only once** to compute the next variable $\mathbf{s}_{k+1} := f_{k+1}(\mathbf{s}_k)$. To compute the VJP of the chain, it then suffices to reverse the order of the operations and to use the corresponding VJPs. Since each function takes a single input, each VJP $\partial f_k(\mathbf{s}_k)^*[\mathbf{r}_k]$ produces a **single** output direction $\mathbf{r}_{k-1} \in \mathcal{S}_k$. A backward pass therefore computes from $k := K$ to $k := 1$, starting from $\mathbf{r}_K := \mathbf{u}$,

$$\mathbf{r}_{k-1} := \partial f_k(\mathbf{s}_{k-1})^*[\mathbf{r}_k] \in \mathcal{S}_{k-1}.$$

The final output is $\mathbf{r}_0 := \partial f(\mathbf{s}_0)^*[\mathbf{u}]$.

In **computation graphs**, intermediate variables may be used more than once, since a node k may have several children nodes. This complicates reversing the computation graph. To circumvent this issue, following the formalism of Roy Frostig (Murphy, 2023, Section 6.2.2.2), we may formally distinguish between the output \mathbf{s}_k of f_k and the inputs $\mathbf{s}_{k \rightarrow j}$ of f_j for $j \in \text{ch}(k)$, the children of f_k . We do so by introducing an operation that simply duplicates \mathbf{s}_k ,

$$\text{dup}(\mathbf{s}_k) := (\mathbf{s}_k, \dots, \mathbf{s}_k).$$

The tuple output by dup is of length $|\text{ch}(k)|$. The forward pass can then be formally rewritten as, for $k \in [K]$,

$$\mathbf{s}_k := f_k(\mathbf{s}_{i_1 \rightarrow k}, \dots, \mathbf{s}_{i_{p_k} \rightarrow k})$$

where $(i_1, \dots, i_{p_k}) := \text{pa}(k)$ followed by

$$(\mathbf{s}_{k \rightarrow j_1}, \dots, \mathbf{s}_{k \rightarrow j_{c_k}}) := \text{dup}(\mathbf{s}_k),$$

where $(j_1, \dots, j_{c_k}) := \text{ch}(k)$. The benefit of this approach is that each duplicated output is input only once to the subsequent child functions. Thanks to this, we can now associate to each variable, $\mathbf{s}_{i \rightarrow j}$ or \mathbf{s}_k , a single corresponding intermediate variable, $\mathbf{r}_{i \rightarrow j}$ or \mathbf{r}_k , in the backward pass. The reverse mode can then simply be written by going through the VJPs of the functions f_k and the VJP of dup in reverse order.

For the VJP of the dup operation, let us denote the intermediate variables in the backward pass associated to $\mathbf{s}_{k \rightarrow j_1}, \dots, \mathbf{s}_{k \rightarrow j_{c_k}}$ by $\mathbf{r}_{k \rightarrow j_1}, \dots, \mathbf{r}_{k \rightarrow j_{c_k}}$. Following the fan-out rule in Proposition 2.9, the VJP of the dup operation on iteration k is

$$\begin{aligned} \mathbf{r}_k &= \text{dup}(\mathbf{s}_k)^*[\mathbf{r}_{k \rightarrow j_1}, \dots, \mathbf{r}_{k \rightarrow j_{c_k}}] \\ &= \sum_{i=1}^{c_k} \text{dup}_i(\mathbf{s}_k)^*[\mathbf{r}_{k \rightarrow j_i}] \\ &= \sum_{i=1}^{c_k} \mathbf{r}_{k \rightarrow j_i}. \end{aligned}$$

In the last line, we used that $\text{dup}_i(\mathbf{s}) = \mathbf{s}$ by definition of the duplication operation so $\partial \text{dup}_i(\mathbf{s}) = \mathbf{I}$ and $\partial \text{dup}_i(\mathbf{s})^* = \mathbf{I}$. The VJP of dup justifies why, if an intermediate value \mathbf{s}_k is used by later functions $f_{j_1}, \dots, f_{j_{c_k}}$, for $(j_1, \dots, j_{c_k}) := \text{ch}(k)$, the derivatives with respect to \mathbf{s}_k need to sum all the variations through the f_j functions into the variable \mathbf{r}_k .

For the VJP of the f_k functions, following the fan-in rule in Proposition 2.8, the VJP returns **multiple** output variations,

$$\mathbf{r}_{i_1 \rightarrow k}, \dots, \mathbf{r}_{i_{p_k} \rightarrow k} = \partial f_k(\mathbf{s}_{i_1 \rightarrow k}, \dots, \mathbf{s}_{i_{p_k} \rightarrow k})^*[\mathbf{r}_k],$$

where $(i_1, \dots, i_{p_k}) := \text{pa}(k)$, and where, for $j \in \{1, \dots, p_k\}$,

$$\mathbf{r}_{i_j \rightarrow k} := \partial_j f_k(\mathbf{s}_{i_1 \rightarrow k}, \dots, \mathbf{s}_{i_{p_k} \rightarrow k})^*[\mathbf{r}_k] \in \mathcal{S}_{i_j}.$$

The overall formalism is summarized in Fig. 8.8.

Implementation

The duplication operation dup is just a formalism to mathematically derive the reverse mode. In practice, the variables \mathbf{s}_k are not duplicated but accessed several times during the backward pass. The variables \mathbf{r}_k are computed by accumulating $\mathbf{r}_{k \rightarrow j}$ into \mathbf{r}_k each time a variable $\mathbf{r}_{k \rightarrow j}$ is computed. Therefore, for each $k \in [K]$, we can compute the VJP and perform the in-place updates,

$$\begin{aligned} \mathbf{r}_{i_1 \rightarrow k}, \dots, \mathbf{r}_{i_{p_k} \rightarrow k} &= \partial f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}})^*[\mathbf{r}_k] \\ \mathbf{r}_{i_j} &\leftarrow \mathbf{r}_{i_j} + \mathbf{r}_{i_j \rightarrow k} \quad \forall j \in \{1, \dots, p_k\}. \end{aligned}$$

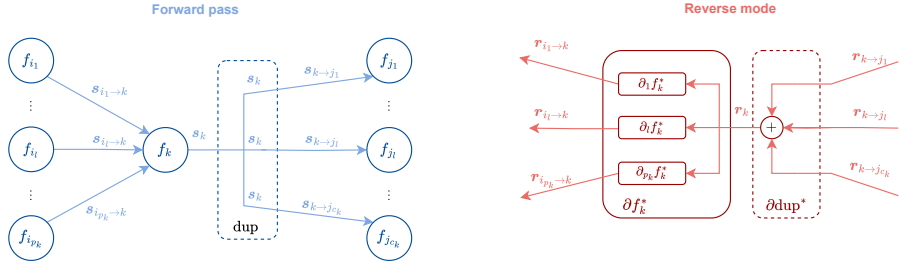


Figure 8.8: Reverse mode automatic differentiation in a computation graph.

Algorithm 8.5 Reverse-mode autodiff for computation graphs

- Functions:** f_1, \dots, f_K in topological order
Inputs: input $s_0 \in \mathcal{S}_0$, output direction $u \in \mathcal{S}_K$
- 1: **for** $k := 1, \dots, K$ **do** ▷ Forward pass
 - 2: Retrieve parent nodes $(i_1, \dots, i_{p_k}) := \text{pa}(k)$
 - 3: Compute $s_k := f_k(s_{i_1}, \dots, s_{i_{p_k}}) \in \mathcal{S}_k$
 - 4: Instantiate VJP $l_k := \partial f_k(s_{i_1}, \dots, s_{i_{p_k}})^*$
 - 5: Initialize $r_K := u$, $r_k \leftarrow 0 \ \forall k \in \{0, \dots, K-1\}$ ▷ Backward pass
 - 6: **for** $k := K, \dots, 1$ **do**
 - 7: Retrieve parent nodes $(i_1, \dots, i_{p_k}) = \text{pa}(k)$
 - 8: Compute $r_{i_1 \rightarrow k}, \dots, r_{i_{p_k} \rightarrow k} := l_k[r_k]$
 - 9: Compute $r_{i_j} \leftarrow r_{i_j} + r_{i_j \rightarrow k} \in \mathcal{S}_{i_j} \ \forall j \in \{1, \dots, p_k\}$
 - 10: **Outputs:** $f(s_0) := s_K \in \mathcal{S}_K$, $\partial f(s_0)^*[u] = r_0 \in \mathcal{S}_0$
-

The topological ordering ensures that r_k has been fully computed when we reach f_k . The resulting generic reverse-mode procedure is presented in Algorithm 8.5.

Example 8.2 (Example of forward and reverse modes). We use Fig. 8.9 to illustrate the forward and reverse modes in a computation graph. Let us assume that the intermediate variables s_1, \dots, s_7 have readily

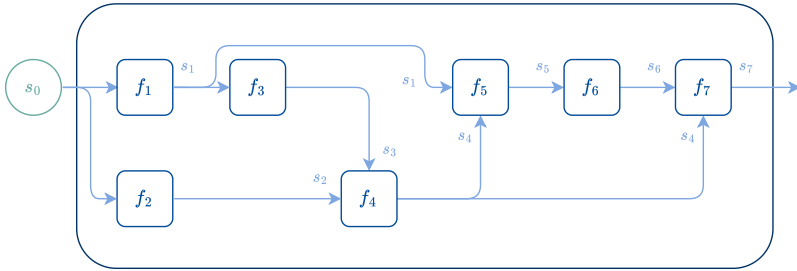


Figure 8.9: Same computation graph as Fig. 4.3 but without the actual definition of each f_k , for simplicity.

been computed. The **forward mode** corresponds to

$$\begin{aligned}
 t_0 &:= v \\
 t_1 &:= \partial f_1(s_0)[t_0] \\
 t_2 &:= \partial f_2(s_0)[t_0] \\
 t_3 &:= \partial f_3(s_1)[t_1] \\
 t_4 &:= \partial f_4(s_2, s_3)[t_2, t_3] = \partial_1 f_4(s_2, s_3)[t_2] + \partial_2 f_4(s_2, s_3)[t_3] \\
 t_5 &:= \partial f_5(s_1, s_4)[t_1, t_4] = \partial_1 f_5(s_1, s_4)[t_1] + \partial_2 f_5(s_1, s_4)[t_4] \\
 t_6 &:= \partial f_6(s_5)[t_5] \\
 t_7 &:= \partial f_7(s_4, s_6)[t_4, t_6] = \partial_1 f_7(s_4, s_6)[t_4] + \partial_2 f_7(s_4, s_6)[t_6].
 \end{aligned}$$

The **reverse mode** corresponds to

$$\begin{aligned}
 (r_{4 \rightarrow 7}, r_{6 \rightarrow 7}) &:= \partial f_7(s_4, s_6)^*[r_7] & r_7 &:= u \\
 r_{5 \rightarrow 6} &:= \partial f_6(s_5)^*[r_6] & r_6 &:= r_{6 \rightarrow 7} \\
 (r_{1 \rightarrow 5}, r_{4 \rightarrow 5}) &:= \partial f_5(s_1, s_4)^*[r_5] & r_5 &:= r_{5 \rightarrow 6} \\
 (r_{2 \rightarrow 4}, r_{3 \rightarrow 4}) &:= \partial f_4(s_2, s_3)^*[r_4] & r_4 &:= r_{4 \rightarrow 5} + r_{4 \rightarrow 7} \\
 r_{1 \rightarrow 3} &:= \partial f_3(s_1)^*[r_3] & r_3 &:= r_{3 \rightarrow 4} \\
 r_{0 \rightarrow 2} &:= \partial f_2(s_0)^*[r_2] & r_2 &:= r_{2 \rightarrow 4} \\
 r_{0 \rightarrow 1} &:= \partial f_1(s_0)^*[r_1] & r_1 &:= r_{1 \rightarrow 3} + r_{1 \rightarrow 5} \\
 & & r_0 &:= r_{0 \rightarrow 1} + r_{0 \rightarrow 2}.
 \end{aligned}$$

8.3.3 Complexity, the Baur-Strassen theorem

For computing the gradient of a function $f: \mathcal{E} \rightarrow \mathbb{R}$ represented by a computation graph, we saw that the reverse mode is more efficient than the forward mode. As we previously stated, assuming that the elementary functions f_k in the DAG and their VJP have roughly the same computational complexity, then f and ∇f have roughly the same computational complexity. This fact is crucial and is the pillar on which modern machine learning relies: it allows us to optimize high-dimensional functions by gradient descent.

For arithmetic circuits, reviewed in Section 4.1.4, this crucial fact is made more precise in the celebrated Baur-Strassen theorem (Baur and Strassen, 1983). If f is a polynomial, then so is its gradient ∇f . The theorem gives an upper bound on the size of the best circuit for computing ∇f from the size of the best circuit for computing f .

Proposition 8.1 (Baur-Strassen’s theorem). For any polynomial $f: \mathcal{E} \rightarrow \mathbb{R}$, we have

$$S(\nabla f) \leq 5 \cdot S(f),$$

where the size $S(f)$ of a polynomial f is defined in Definition 4.1.

A simpler proof by backward induction was given by Morgenstern (1985). See also the proof of Theorem 9.10 in Chen *et al.* (2011). For general computation graphs, that have more primitive functions than just $+$ and \times , a similar result can be obtained; see, e.g., (Bolte *et al.*, 2022, Theorem 2).

8.4 Implementation

8.4.1 Primitive functions

An autodiff system implements a set \mathcal{A} of primitive or elementary functions, which serve as building blocks for creating other functions, by function composition. For instance, we saw that in arithmetic circuits (Section 4.1.4), $\mathcal{A} = \{+, \times\}$. More generally, \mathcal{A} may contain all the necessary functions for expressing programs. We emphasize, however, that \mathcal{A} is not necessarily restricted to low-level functions such as \log

and \exp , but may also contain higher-level functions. For instance, even though the \log - \sum - \exp can be expressed as the composition of elementary operations (\log , \sum , \exp), it is usually included as a primitive on its own, both because it is a very commonly-used building block, but also for numerical stability reasons.

8.4.2 Closure under function composition

Each function f_k in a computation graph belongs to a set \mathcal{F} , the class of functions supported by the system. A desirable property of an autodiff implementation is that the set \mathcal{F} is **closed** under function composition, meaning that if $f \in \mathcal{F}$ and $g \in \mathcal{F}$, then $f \circ g \in \mathcal{F}$. This means that composed functions can themselves be used for composing new functions. This property is also crucial for supporting higher-order differentiation (Chapter 9) and automatic linear transposition (Section 8.4.4). When f_k is a composition of elementary functions in \mathcal{A} , then f_k itself is a **nested** DAG. However, we can always **inline** each composite function, such that all functions in the DAG belong to \mathcal{A} .

8.4.3 Examples of JVPs and VJPs

An autodiff system must implement for each $f \in \mathcal{A}$ its JVP for supporting the forward mode, and its VJP for supporting the reverse mode. We give a couple of examples. We start with the JVP and VJP of linear functions.

Example 8.3 (JVP and VJP of linear functions). Consider the matrix-vector product $f(\mathbf{W}) = \mathbf{W}\mathbf{x} \in \mathbb{R}^M$, where $\mathbf{x} \in \mathbb{R}^D$ is fixed and $\mathbf{W} \in \mathbb{R}^{M \times D}$. As already mentioned in Section 2.3.1, the JVP of f at $\mathbf{W} \in \mathbb{R}^{M \times D}$ along an input direction $\mathbf{V} \in \mathbb{R}^{M \times D}$ is simply

$$\partial f(\mathbf{W})[\mathbf{V}] = f(\mathbf{V}) = \mathbf{V}\mathbf{x} \in \mathbb{R}^M.$$

To find the associated VJP, we note that for any $\mathbf{u} \in \mathbb{R}^M$ and $\mathbf{V} \in \mathbb{R}^{M \times D}$, we must have $\langle \partial f(\mathbf{W})[\mathbf{V}], \mathbf{u} \rangle = \langle \mathbf{V}, \partial f(\mathbf{W})^*[\mathbf{u}] \rangle$. Using the properties of the trace, we have

$$\langle \partial f(\mathbf{W})[\mathbf{V}], \mathbf{u} \rangle = \langle \mathbf{V}\mathbf{x}, \mathbf{u} \rangle = \text{tr}(\mathbf{x}^\top \mathbf{V}^\top \mathbf{u}) = \langle \mathbf{V}, \mathbf{u}\mathbf{x}^\top \rangle.$$

Therefore, we find that the VJP is given by

$$\partial f(\mathbf{W})^*[\mathbf{u}] = \mathbf{u}\mathbf{x}^\top \in \mathbb{R}^{M \times D}.$$

Similarly, consider now a matrix-matrix product $f(\mathbf{W}) = \mathbf{W}\mathbf{X}$, where $\mathbf{W} \in \mathbb{R}^{M \times D}$ and where $\mathbf{X} \in \mathbb{R}^{D \times N}$ is fixed. The JVP at $\mathbf{W} \in \mathbb{R}^{M \times D}$ along an input direction $\mathbf{V} \in \mathbb{R}^{M \times D}$ is simply

$$\partial f(\mathbf{W})[\mathbf{V}] = f(\mathbf{V}) = \mathbf{V}\mathbf{X} \in \mathbb{R}^{M \times N}.$$

The VJP along the output direction $\mathbf{U} \in \mathbb{R}^{M \times N}$ is

$$\partial f(\mathbf{W})^*[\mathbf{U}] = \mathbf{U}\mathbf{X}^\top \in \mathbb{R}^{M \times D}.$$

Another simple example are element-wise separable functions.

Example 8.4 (JVP and VJP of separable function). Consider the function $f(\mathbf{w}) := (g_1(w_1), \dots, g_P(w_P))$, where each $g_i: \mathbb{R} \rightarrow \mathbb{R}$ has a derivative g'_i . The Jacobian matrix is then a diagonal matrix

$$\partial f(\mathbf{w}) = \mathbf{diag}(g'_1(w_1), \dots, g'_P(w_P)) \in \mathbb{R}^{P \times P}.$$

In this case, the JVP and VJP are actually the same

$$\partial f(\mathbf{w})[\mathbf{v}] = \partial f(\mathbf{w})^*[\mathbf{v}] = (g'_1(w_1), \dots, g'_P(w_P)) \odot \mathbf{v},$$

where \odot indicates element-wise multiplication.

8.4.4 Automatic linear transposition

On first sight, if we want to support both forward and reverse modes, it appears like we need to implement both the JVP and the VJP for each primitive operation $f \in \mathcal{A}$. Fortunately, there exists a way to recover VJPs from JVPs, and vice-versa.

We saw in Section 2.3 that if $l(\mathbf{w})$ is a linear map, then its JVP is $\partial l(\mathbf{w})[\mathbf{v}] = l(\mathbf{v})$ (independent of \mathbf{w}). Conversely, the VJP is $\partial l(\mathbf{w})^*[\mathbf{u}] = l^*(\mathbf{u})$, where l^* is the adjoint operator of l (again, independent of \mathbf{w}).

Let us define $l(\mathbf{u}; \mathbf{w}) := \partial f(\mathbf{w})^*[\mathbf{u}]$, i.e., the VJP of f in the output direction \mathbf{u} . Since $l(\mathbf{u}; \mathbf{w})$ is linear in \mathbf{u} , we can apply the reasoning

above to compute its VJP

$$\partial l(\mathbf{u}; \mathbf{w})^*[\mathbf{v}] = l^*(\mathbf{v}; \mathbf{w}) = \partial f(\mathbf{w})^{**}[\mathbf{v}] = \partial f(\mathbf{w})[\mathbf{v}],$$

which is independent of \mathbf{u} . In words, the VJP of a VJP is the corresponding JVP! This means that we can implement forward-mode autodiff even if we only have access to VJPs. As an illustration and sanity check, we give the following example.

Example 8.5 (Automatic transpose of “dot”). If we define $f(\mathbf{x}, \mathbf{W}) := \mathbf{W}\mathbf{x}$, from Example 8.3, we know that

$$\begin{aligned} \partial f(\mathbf{x}, \mathbf{W})^*[\mathbf{u}] &= (\mathbf{W}^\top \mathbf{u}, \mathbf{u}\mathbf{x}^\top) \\ &= (f(\mathbf{u}, \mathbf{W}^\top), f(\mathbf{x}^\top, \mathbf{u})) \\ &=: l(\mathbf{u}; \mathbf{x}, \mathbf{W}). \end{aligned}$$

Using Proposition 2.9, we obtain

$$\begin{aligned} \partial l(\mathbf{u}; \mathbf{x}, \mathbf{W})^*[\mathbf{v}, \mathbf{V}] &= f(\mathbf{v}, \mathbf{W}) + f(\mathbf{x}, \mathbf{V}) \\ &= \mathbf{W}\mathbf{v} + \mathbf{V}\mathbf{x} \\ &= \partial f(\mathbf{x}, \mathbf{W})[\mathbf{v}, \mathbf{V}]. \end{aligned}$$

The other direction, automatically creating a VJP from a JVP, is also possible but is more technical and relies on the notion of **partial evaluation** (Frostig *et al.*, 2021; Radul *et al.*, 2022).

8.5 Checkpointing

We saw that forward-mode autodiff can release intermediate computations from memory along the way, while reverse-mode autodiff needs to cache all of them. This means that the memory complexity of reverse-mode autodiff, in its standard form, grows linearly with the number of nodes in the computation graph. A commonly-used technique to circumvent this issue is checkpointing, which trades-off computation time for better memory usage. Checkpointing works by selectively storing only a subset of the intermediate values, called **checkpoints**, and by recomputing others on the fly. The specific choice of the checkpoint locations in the computation graph determines the memory-computation

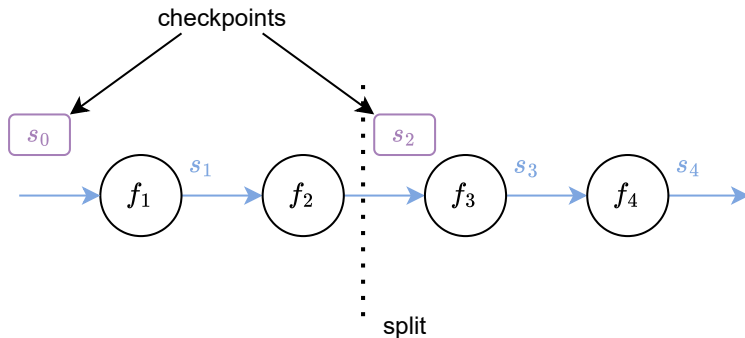


Figure 8.10: Checkpointing trades-off computation time for better memory usage by selectively storing only a subset of the intermediate values, called checkpoints, and by recomputing others on the fly. Recursive halving and dynamic programming are two divide-and-conquer strategies to select the checkpoint locations.

trade-off. While it is possible to heuristically set checkpoints at user-specified locations, it is also possible to perform a checkpointing strategy algorithmically, as studied in depth by Griewank (1992) and Griewank and Walther (2008). In this section, we review two divide-and-conquer algorithms: recursive halving and dynamic programming. Our exposition focuses on computation chains $f = f_K \circ \dots \circ f_1$, with $f_k : \mathbb{R}^D \rightarrow \mathbb{R}^D$ for simplicity.

Computational and memory complexities at two extremes. Let $\mathcal{C}(K)$ be the number of calls to the individual functions f_k (we ignore the cost of computing the intermediate VJPs) and $\mathcal{M}(K)$ be the number of function inputs cached, when performing reverse-mode autodiff on a chain $f = f_K \circ \dots \circ f_1$. On one extreme, if we store all intermediate computations, as done in Algorithm 8.1, to compute only the VJP $\partial f(s_0)^*[u]$, we have

$$\mathcal{C}(K) = K - 1 \quad \text{and} \quad \mathcal{M}(K) = K.$$

This is optimal w.r.t. computational complexity, but suboptimal w.r.t. memory. On the other extreme, if we only store the initial input, as

Algorithm 8.6 Reverse-mode autodiff with constant memory

$\text{vjp_full_recompute}(f_K \circ \dots \circ f_1, \mathbf{s}_0, \mathbf{u}) := \partial(f_K \circ \dots \circ f_1)(\mathbf{s}_0)^*[\mathbf{u}]$

Inputs: Chain $f_K \circ \dots \circ f_1$, input $\mathbf{s}_0 \in \mathcal{S}_0$, output direction $\mathbf{u} \in \mathcal{S}_K$

- 1: **if** $K = 1$ **then**
- 2: **return** $\partial f_1(\mathbf{s}_0)^*[\mathbf{u}]$
- 3: **else**
- 4: Set $\mathbf{r}_K = \mathbf{u}$
- 5: **for** $k := K, \dots, 1$ **do**
- 6: Compute $\mathbf{s}_{k-1} = (f_{k-1} \circ \dots \circ f_1)(\mathbf{s}_0)$
- 7: Compute $\mathbf{r}_{k-1} = \partial f_k(\mathbf{s}_{k-1})^*[\mathbf{r}_k]$
- 8: **return:** \mathbf{r}_0

done in Algorithm 8.6, then we have

$$\mathcal{C}(K) = K(K-1)/2 \quad \text{and} \quad \mathcal{M}(K) = 1.$$

This is optimal w.r.t. memory but leads to a computational complexity that is quadratic in K .

8.5.1 Recursive halving

As a first step towards obtaining a better computation-memory trade-off, we may split the chain $\mathbf{s}_K = f_K \circ \dots \circ f_1(\mathbf{s}_0)$ as

$$\begin{aligned} \mathbf{s}_{K/2} &= f_{K/2} \circ \dots \circ f_1(\mathbf{s}_0) \\ \mathbf{s}_K &= f_K \circ \dots \circ f_{K/2+1}(\mathbf{s}_{K/2}), \end{aligned}$$

for K even. Then, rather than recomputing all intermediate computations \mathbf{s}_k from the input \mathbf{s}_0 as in Algorithm 8.6, we can store $\mathbf{s}_{K/2}$ and recompute \mathbf{s}_k for $k > K/2$ starting from $\mathbf{s}_{K/2}$. Formally, this strategy amounts to the following steps.

1. Compute $\mathbf{s}_{K/2} = f_{K/2} \circ \dots \circ f_1(\mathbf{s}_0)$
2. Compute $\mathbf{r}_{K/2} = \text{vjp_full_recompute}(f_K \circ \dots \circ f_{K/2+1}, \mathbf{s}_{K/2}, \mathbf{u})$
3. Compute $\mathbf{r}_0 = \text{vjp_full_recompute}(f_{K/2} \circ \dots \circ f_1, \mathbf{s}_0, \mathbf{r}_{K/2})$

Algorithm 8.7 Reverse-mode autodiff with recursive halving

$\text{vjp_halving}(f_K \circ \dots \circ f_1, \mathbf{s}_0, \mathbf{u}) := \partial(f_K \circ \dots \circ f_1)(\mathbf{s}_0)^*[\mathbf{u}]$

Functions: Chain $f_K \circ \dots \circ f_1$

Inputs: input $\mathbf{s}_0 \in \mathcal{S}_0$, output direction $\mathbf{u} \in \mathcal{S}_K$

```

1: if  $K = 1$  then
2:   return  $\partial f_1(\mathbf{s}_0)^*[\mathbf{u}]$ 
3: else
4:   Compute  $\mathbf{s}_{K/2} = f_{K/2} \circ \dots \circ f_1(\mathbf{s}_0)$ 
5:   Compute  $\mathbf{r}_{K/2} = \text{vjp\_halving}(f_K \circ \dots \circ f_{K/2+1}, \mathbf{s}_{K/2}, \mathbf{u})$ 
6:   Compute  $\mathbf{r}_0 = \text{vjp\_halving}(f_{K/2} \circ \dots \circ f_1, \mathbf{s}_0, \mathbf{r}_{K/2})$ 
7:   return:  $\mathbf{r}_0$ 

```

At the expense of having to store the additional checkpoint $\mathbf{s}_{K/2}$, this already roughly halves the computational complexity compared to Algorithm 8.6.

We can then apply this reasoning recursively, as formalized in Algorithm 8.7. The algorithm is known as **recursive binary schedule** (Griewank, 2003) and illustrated in Fig. 8.11. In terms of number of function evaluations $\mathcal{C}(K)$, for K even, we make $K/2$ function calls, and we call the procedure recursively twice, that is,

$$\mathcal{C}(K) = 2\mathcal{C}(K/2) + K/2.$$

If the chain is of length 1, we directly use the VJP, so $\mathcal{C}(1) = 0$. Hence, the numbers of function calls, if K is a power of 2, is

$$\mathcal{C}(K) = \frac{K}{2} \log_2 K.$$

In terms of memory usage, Algorithm 8.7 uses \mathbf{s}_0 not only at line 4 but also at line 6. So when the algorithm is called recursively on the second half of the chain at line 5, one memory slot is taken by \mathbf{s}_0 . This line is called recursively until the chain is reduced to a single function. At that point, the total number of memory slots used is equal to the number of times we split the function in half, that is, $\log_2 K$ for K a power of 2. On the other hand, the input \mathbf{s}_0 is no longer used after line 6 of Algorithm 8.7. At that line, the memory slot taken by \mathbf{s}_0 can be consumed by the recursive call on the first-half. In other words, calling

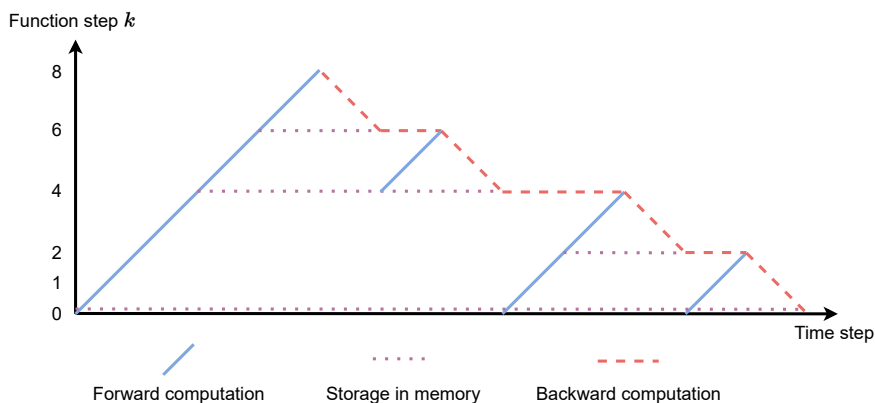


Figure 8.11: Illustration of checkpointing with recursive halving, for a chain of 8 functions. The chain is first fully evaluated while storing some computations as checkpoints in memory. Then, during the backward pass, we recompute some intermediate values from the latest checkpoint available. In contrast, vanilla reverse-mode autodiff (with full caching of the intermediate computations) would lead to a simple triangle shape.

the algorithm recursively on the first half does not incur extra memory cost. So if K is a power of 2, the memory cost of Algorithm 8.7 is

$$\mathcal{M}(K) = \log_2 K.$$

8.5.2 Dynamic programming

Recursive halving requires $\log_2 K$ memory slots for a chain of length K . However, as illustrated in Fig. 8.11, at a given time step, all memory slots may not be exploited.

To optimize the approach, we observe that recursive halving is just one instance of a program that splits the chain and calls itself recursively on each part. In other words, it is a form of **divide-and-conquer** algorithm. Rather than splitting the chain in half, we may consider splitting the chain at some index l . One split is used to reverse the computations from $l + 1$ to K by a recursive call that consumes one memory slot. The other split is used on a recursive call that reverses the computations from 0 to l . That second call does not require an additional memory slot, as it can use directly the memory slot used

by the original input s_0 . To split the chain in such two parts, we need l intermediate computations to go from s_0 to s_l . The computational complexity $\mathcal{C}(k, s)$, counted as the number of function evaluations, for a chain of length k with s memory slots then satisfies the recurrence

$$\mathcal{C}(k, s) = \mathcal{C}(k - l, s - 1) + \mathcal{C}(l, s) + l,$$

for all $l \in \{1, \dots, k - 1\}$. By simply taking $l = k/2$, we recover exactly the computational complexity of recursive halving. To refine the latter, we may split the chain by selecting l to minimize the complexity. An optimal scheme must satisfy the recursive equation,

$$\mathcal{C}^*(k, s) := \min_{1 \leq l \leq K-1} \{\mathcal{C}^*(k - l, s - 1) + \mathcal{C}^*(l, s) + l\}. \quad (8.3)$$

Note that $\mathcal{C}^*(K, S)$ can be computed from $\mathcal{C}^*(k, s)$ for $k = 1, \dots, K - 1$, $s = 1, \dots, S - 1$. This suggests a **dynamic programming** approach to find an optimal scheme algorithmically. For a chain of length $k = 1$, the cost is null as we directly reverse the computation, so $\mathcal{C}^*(1, s) := 0$. On the other hand for a memory $s = 1$, there is only one possible scheme that saves only the initial input as in Algorithm 8.6, so $\mathcal{C}^*(k, 1) := (k(k - 1))/2$. The values $\mathcal{C}^*(k, s)$ can then be computed incrementally from $k = 1$ to K and $s = 1$ to S using Eq. (8.3). The optimal splits can be recorded along the way as

$$l^*(k, s) := \arg \min_{1 \leq l \leq k-1} \{\mathcal{C}^*(k - l, s - 1) + \mathcal{C}^*(l, s) + l\}.$$

The optimal split for K, S can then be found by **backtracking** the optimal splits along both branches corresponding to $\mathcal{C}^*(k - l, s - 1)$ and $\mathcal{C}^*(l, s)$. As the final output consists in traversing a binary tree, it was called **treeverse** (Griewank, 1992). Note that the dynamic programming procedure is generic and could a priori incorporate varying computational costs for the intermediate functions f_k .

Analytical formula

It turns out that we can also find an optimal scheme *analytically*. This scheme was found by Griewank (1992), following the analysis of optimal inversions of sequential programs by divide-and-conquer algorithms

done by Grimm *et al.* (1996); see also Griewank (2003, Section 6) for a simple proof. The main idea consists in considering the number of times an evaluation step f_k is repeated. As we split the chain at l , all steps from 1 to l will be repeated at least once. In other words, treating the second half of the chain incurs one memory cost, while treating the first half of the chain incurs one repetition cost. Griewank (1992) shows that for fixed K, S , we can find the minimal number of repetitions analytically and build the corresponding scheme with simple formulas for the optimal splits.

Compared to the dynamic programming approach, it means that we do not need to compute the pointers $l^*(k, s)$, and we can use a simple formula to set $l^*(k, s)$. We still need to traverse the corresponding binary tree given K, S and $l^*(k, s)$ to obtain the schedules. Note that such optimal scheme does not take into account varying computational costs for the functions f_k .

8.5.3 Online checkpointing

The optimal scheme presented above requires knowing the total number of nodes in the computation graph ahead of time. However, when differentiating through for example a while loop (Section 5.10), this is not the case. To circumvent this issue, online checkpointing schemes have been developed and proven to be nearly optimal (Stumm and Walther, 2010; Wang *et al.*, 2009). These schemes start by defining a set of S checkpoints with the first S computations, then these checkpoints are rewritten dynamically as the computations keep going. Once the computations terminate, the optimal approach presented above for a fixed length is applied on the set of checkpoints recorded.

8.6 Reversible layers

8.6.1 General case

The memory requirements of reverse-mode autodiff can be completely alleviated when the functions f_k are invertible (meaning that f_k^{-1} exists) and when f_k^{-1} is easily accessible. In that case, rather than storing the intermediate computations \mathbf{s}_{k-1} , necessary to compute the VJP $\mathbf{r}_k \mapsto$

Algorithm 8.8 Reverse-mode autodiff for reversible chains.

Functions: $f := f_K \circ \dots \circ f_1$, with each f_k invertible

Inputs: input $\mathbf{s}_0 \in \mathcal{S}_0$, output direction $\mathbf{u} \in \mathcal{S}_K$

- 1: Compute $\mathbf{s}_K = f_K \circ \dots \circ f_1(\mathbf{s}_0)$
- 2: **for** $k := K, \dots, 1$ **do**
- 3: Compute $\mathbf{s}_{k-1} = f_k^{-1}(\mathbf{s}_k)$
- 4: Compute $\mathbf{r}_{k-1} = \partial f_k(\mathbf{s}_{k-1})^*[\mathbf{r}_k]$

Outputs: $f(\mathbf{s}_0) := \mathbf{s}_K$, $\partial f(\mathbf{s}_0)^*[\mathbf{u}] = \mathbf{r}_0$

$\partial f_k(\mathbf{s}_{k-1})^*[\mathbf{r}_k]$, one can compute them on the fly during the backward pass from \mathbf{s}_k using $\mathbf{s}_{k-1} = f_k^{-1}(\mathbf{s}_k)$. We summarize the procedure for the case of computation chains in Algorithm 8.8. Compared to vanilla reverse-mode autodiff in Algorithm 8.2, the algorithm has optimal memory complexity, as we can release \mathbf{s}_k and \mathbf{r}_k as we go.

In practice, f_k^{-1} often does not exist or may not be easily accessible. However, network architectures can be constructed to be easily invertible by design. Examples include reversible residual networks (Gomez *et al.*, 2017), orthonormal RNNs (Helfrich *et al.*, 2018), neural ODEs (Section 12.6), and momentum residual neural networks (Sander *et al.*, 2021a); see also references therein.

8.6.2 Case of orthonormal JVPs

When the JVP of each f_k is an orthonormal linear mapping, i.e.,

$$\partial f_k(\mathbf{s}_{k-1})^{-1} = \partial f_k(\mathbf{s}_{k-1})^*,$$

it is easy to check that the VJP of $f = f_K \circ \dots \circ f_1$ is equal to the JVP of $f^{-1} = f_1^{-1} \circ \dots \circ f_K^{-1}$, that is

$$\partial f(\mathbf{s}_0)^*[\mathbf{u}] = \partial f^{-1}(\mathbf{s}_K)[\mathbf{u}].$$

In other words, in the case of orthonormal JVPs, reverse-mode autodiff of f coincides with forward-mode autodiff of f^{-1} .

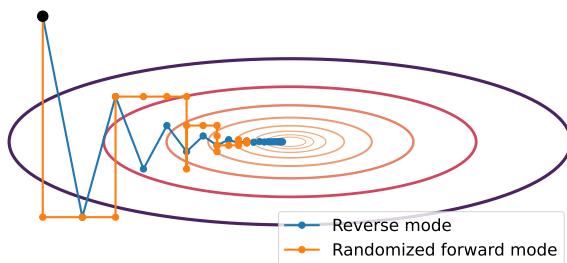


Figure 8.12: The randomized forward-mode gradient estimator only requires forward passes, but it suffers from high variance, even more so in high dimension.

8.7 Randomized forward-mode gradient estimator

Forward-mode autodiff does not require to store intermediate activations. However, for a function $f: \mathbb{R}^P \rightarrow \mathbb{R}$, computing the gradient ∇f using forward-mode autodiff requires P JVPs, which is intractable if P is large. Can we approximate ∇f with fewer JVPs? The following proposition gives an **unbiased estimator** of ∇f that only involves JVPs.

Proposition 8.2 (Unbiased forward-mode estimator of the gradient).

Let $f: \mathbb{R}^P \rightarrow \mathbb{R}$ be a differentiable function. Then,

$$\begin{aligned}\nabla f(\boldsymbol{\mu}) &= \mathbb{E}_{Z \sim p} [\partial f(\boldsymbol{\mu})[Z]Z] \\ &= \mathbb{E}_{Z \sim p} [\langle \nabla f(\boldsymbol{\mu}), Z \rangle Z].\end{aligned}$$

where $p := \text{Normal}(0, 1)^P$ is the isotropic Gaussian distribution.

This estimator is for instance used by Baydin *et al.* (2022). It can be seen as the **zero-temperature limit** of the gradient of a perturbed function, estimated by the score-function estimator (SFE); see Section 14.4.6.

In practice, the expectation above can be approximated by drawing M noise vectors $\mathbf{z}_1, \dots, \mathbf{z}_M$, and averaging $\langle \nabla f(\boldsymbol{\mu}), \mathbf{z}_i \rangle$ over $i \in [M]$.

A word of caution: while this estimator can be useful for example when we do not want to store the intermediate activations for memory reasons, this of course comes at the cost of increasing the variance, which influences the convergence rate of SGD, as seen in Section 16.2.

8.8 Summary

- Computer programs can be seen as directed acyclic graphs, where nodes correspond to the output of intermediate operations in the program, and edges represent the dependencies of current operations on past operations.
- Automatic differentiation (autodiff) for a function $f: \mathbb{R}^P \rightarrow \mathbb{R}^M$ has two main modes: forward mode and reverse mode.
- The forward mode: i) uses JVPs, ii) builds the Jacobian one column at a time, iii) is efficient for tall Jacobians ($M \geq P$), iv) need not store intermediate computations.
- The reverse mode: i) uses VJPs, builds the Jacobian one row at a time, iii) is efficient for wide Jacobians ($P \geq M$), iv) needs to store intermediate computations, in order to be computationally optimal.
- To trade computational efficiency for better memory efficiency, we can use checkpointing techniques.
- The complexity of computing the gradient of a function $f: \mathbb{R}^P \rightarrow \mathbb{R}$ using the reverse mode is at most a constant time bigger than that of evaluating the function itself. This is the Baur-Strassen theorem, in arithmetic circuits. This astonishing result is one of the pillars of modern machine learning.

9

Second-order automatic differentiation

We review in this chapter how to perform automatic differentiation for second-order derivatives.

9.1 Hessian-vector products

We consider in this section a function $f: \mathcal{E} \rightarrow \mathbb{R}$. Similarly to the Jacobian, for most purposes, we do not need access to the full Hessian but rather to the Hessian-vector product (HVP) $\nabla^2 f(\mathbf{w})[\mathbf{v}]$ at $\mathbf{w} \in \mathcal{E}$, in a direction $\mathbf{v} \in \mathcal{E}$, as defined in Definition 2.19. The latter can be computed in four different ways, depending on how we combine the two main modes of autodiff.

9.1.1 Four possible methods

An HVP can be computed in four different ways.

1. **Reverse on reverse:** The Hessian can be seen as the transposed Jacobian of the gradient, hence the HVP can be computed as the **VJP of the gradient**,

$$\nabla^2 f(\mathbf{w})[\mathbf{v}] = \partial(\nabla f)(\mathbf{w})^*[\mathbf{v}].$$

2. **Forward on reverse:** Owing to its symmetry (see Proposition 2.10), the Hessian can also be seen as the Jacobian of the gradient, hence the HVP can be computed as the **JVP of the gradient**,

$$\nabla^2 f(\mathbf{w})[\mathbf{v}] = \partial(\nabla f)(\mathbf{w})[\mathbf{v}].$$

3. **Reverse on forward:** Recall that for any function $g: \mathcal{E} \rightarrow \mathcal{E}$, the VJP can equivalently be defined as the gradient along an output direction $\mathbf{v} \in \mathcal{E}$, that is,

$$\partial g(\mathbf{w})^*[\mathbf{v}] = \nabla \langle g, \mathbf{v} \rangle(\mathbf{w}),$$

where we recall the shorthand $\langle g, \mathbf{v} \rangle(\mathbf{w}) := \langle \mathbf{v}, g(\mathbf{w}) \rangle$, so that $\langle g, \mathbf{v} \rangle$ is a function of \mathbf{w} . In our case, we can therefore rewrite the reverse-on-reverse approach as

$$\partial(\nabla f)(\mathbf{w})^*[\mathbf{v}] = \nabla \langle \nabla f, \mathbf{v} \rangle(\mathbf{w}).$$

We know that $\langle \nabla f, \mathbf{v} \rangle(\mathbf{w}) = \langle \nabla f(\mathbf{w}), \mathbf{v} \rangle = \partial f(\mathbf{w})[\mathbf{v}]$ is the JVP of f at \mathbf{w} along \mathbf{v} . Therefore, we can also compute the HVP as the **gradient of the JVP** of f at \mathbf{w} along \mathbf{v} ,

$$\nabla^2 f(\mathbf{w})[\mathbf{v}] = \nabla(\partial f(\cdot)[\mathbf{v}])(\mathbf{w}),$$

where we use the notation $(\partial f(\cdot)[\mathbf{v}])(\mathbf{w}) := \partial f(\mathbf{w})[\mathbf{v}]$ to insist on the fact that it is a function of \mathbf{w} .

4. **Forward on forward:** Finally, we can use the definition of the HVP in Definition 2.19 as a vector of second partial derivatives along \mathbf{v} and each canonical direction. That is, assuming $\mathcal{E} = \mathbb{R}^P$, we can compute the **JVP of the JVP** P times,

$$\nabla^2 f(\mathbf{w})[\mathbf{v}] = (\partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{e}_i])_{i=1}^P.$$

The four different ways of computing the HVP are summarized in Table 9.1.

9.1.2 Complexity

To get a sense of the computational and memory complexity of the four approaches, we consider a chain of functions $f := f_K \circ \cdots \circ f_1$ as done

Method	Computation
Reverse on reverse (VJP of gradient)	$\partial(\nabla f)(\mathbf{w})^*[\mathbf{v}]$
Forward on reverse (JVP of gradient)	$\partial(\nabla f)(\mathbf{w})[\mathbf{v}]$
Reverse on forward (gradient of JVP)	$\nabla(\partial f(\cdot)[\mathbf{v}])(\mathbf{w})$
Forward on forward (JVPs of JVPs)	$(\partial^2 f(\mathbf{w})[\mathbf{v}, \mathbf{e}_i])_{i=1}^P$

Table 9.1: Four different ways of computing the HVP $\nabla^2 f(\mathbf{w})[\mathbf{v}]$.

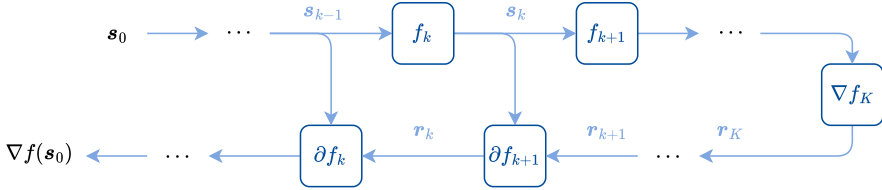


Figure 9.1: Computation graph corresponding to reverse mode autodiff for evaluating the gradient of $f = f_K \circ \dots \circ f_1$. While f is a simple chain, ∇f is a DAG.

in Section 8.1. To simplify our analysis, we assume $f_k: \mathbb{R}^P \rightarrow \mathbb{R}^P$ for $k \in \{1, \dots, K-1\}$ and $f_K: \mathbb{R}^P \rightarrow \mathbb{R}$.

The computation graph of the reverse mode is illustrated in Fig. 9.1. While $f = f_K \circ \dots \circ f_1$ would be represented by a simple chain, the computational graph of ∇f is no longer a chain: it is a DAG. This is due to the computations of $\partial f_k(\mathbf{s}_{k-1})[\mathbf{r}_k]$, where both \mathbf{s}_{k-1} and \mathbf{r}_k depend on \mathbf{s}_0 .

We illustrate the computation graphs of reverse-on-reverse and forward-on-reverse in Fig. 9.2 and Fig. 9.3 respectively. By applying reverse mode on reverse mode, at each fan-in operation $\mathbf{s}_{k-1}, \mathbf{r}_k \mapsto \partial f_k(\mathbf{s}_{k-1})[\mathbf{r}_k]$, the reverse mode on ∇f branches out in two paths that are later merged by a sum. By applying forward mode on top of reverse mode, the flow of computations simply follows the one of ∇f .

With this in mind, following a similar calculation as for Table 8.1, we obtain the following results. We assume that each $\partial f_k(\mathbf{s}_{k-1})$ is a dense linear operator, so that its application has the same cost as a matrix-vector multiplication. For the memory complexity, we consider that the inputs of each operation is saved to compute the required

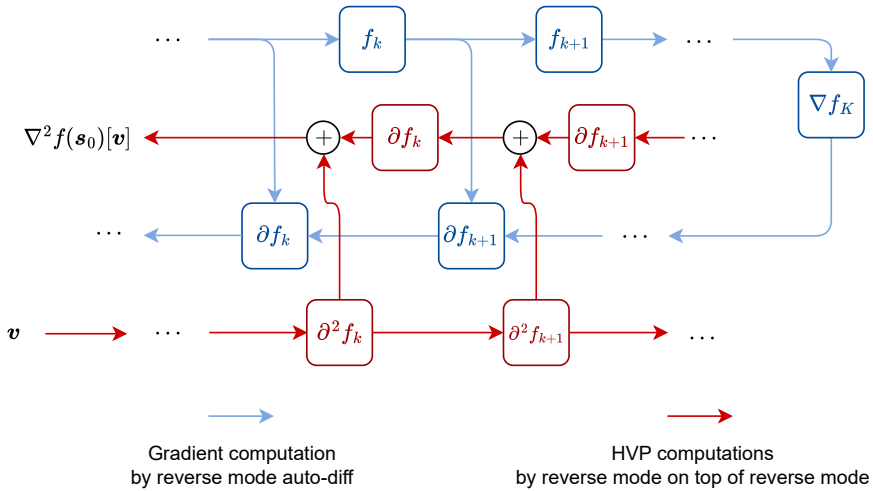


Figure 9.2: Computation graph for computing the HVP $\nabla^2 f(x)[v]$ by using reverse mode on top of reverse mode. As the computation graph of ∇f induces fan-in operations $s_{k-1}, r_k \mapsto \partial f_k(s_{k-1})[r_k]$, the reverse mode applied on ∇f induces branching of the computations at each such node.

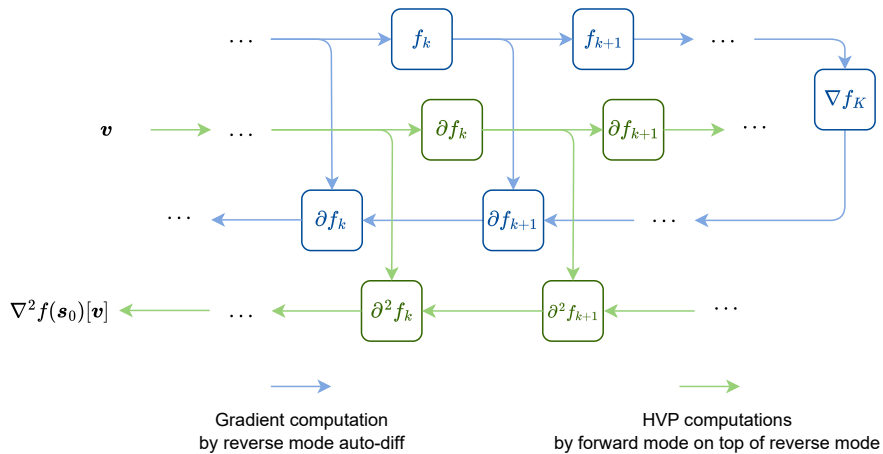


Figure 9.3: Computation graph for computing the HVP $\nabla^2 f(x)[v]$ by using forward mode on top of reverse mode. The forward mode naturally follows the computations done for the gradient, except that it passes through the derivatives of the intermediate operations.

derivatives in the backward passes.

1. **Reverse on reverse:** $O(KP^2)$ time and $O(KP)$ space.
2. **Forward on reverse:** $O(KP^2)$ time and $O(KP)$ space.
3. **Reverse on forward:** $O(KP^2)$ time and $O(KP)$ space.
4. **Forward on forward:** $O(KP^3)$ time and $O(3P)$ space for the P JVPs with e_1, \dots, e_P .

We see that, for chains of functions, “reverse on reverse”, “forward on reverse” and “reverse on forward” all have similar time complexities up to some constant factors. Using reverse mode on top of reverse mode requires storing the information backpropagated, i.e., the \mathbf{r}_k (resp. the information forwarded, i.e., the \mathbf{t}_k in Fig. 8.1), to perform the final reverse pass. By using forward mode on top of reverse mode, this additional cost is not incurred, making it slightly less memory expensive. In addition, reverse mode on top of reverse mode induces a few additional summations due to the branching and merge operations depicted in Fig. 9.2. The same holds when using reverse on top of forward as we cannot avoid fan-in operations (this time of the form $\mathbf{s}_{k-1}, \mathbf{t}_{k-1} \mapsto \partial f_k(\mathbf{s}_{k-1})[\mathbf{t}_{k-1}]$). Unfortunately, “forward on forward” is prohibitively expensive.

To summarize, among the four approaches presented to compute HVPs, the forward-over-reverse mode is a priori the most preferable in terms of computational and memory complexities. Note, however, that computations of higher derivatives can benefit from dedicated autodiff implementations such as Taylor mode autodiff, that do not merely compose forward and reverse modes. For general functions f , it is reasonable to benchmark the first three methods to determine which method is the best for the function at hand.

9.2 Gauss-Newton matrix

9.2.1 An approximation of the Hessian

The Hessian matrix $\nabla^2 L(\mathbf{w})$ of a function $L: \mathcal{W} \rightarrow \mathbb{R}$ is often used to construct a quadratic approximation of $L(\mathbf{w})$,

$$L(\mathbf{w} + \mathbf{v}) \approx \langle \nabla L(\mathbf{w}), \mathbf{v} \rangle + \frac{1}{2} \langle \mathbf{v}, \nabla^2 L(\mathbf{w}) \mathbf{v} \rangle.$$

Unfortunately, when L is nonconvex, $\nabla^2 L(\mathbf{w})$ is typically an **indefinite matrix**, which means that the above approximation is a **nonconvex quadratic** w.r.t. \mathbf{v} . For instance, if $L = \ell \circ f$ with ℓ convex, then L is convex if f is linear, but it is typically nonconvex if f is nonlinear. The (generalized) Gauss-Newton matrix is a principled alternative to the Hessian, which is defined for $L := \ell \circ f$.

Definition 9.1 (Gauss-Newton matrix). Given a differentiable function $f: \mathcal{W} \rightarrow \mathcal{M}$ and a twice differentiable function $\ell: \mathcal{M} \rightarrow \mathbb{R}$, the (generalized) Gauss-Newton matrix of the composition $L = \ell \circ f$ evaluated at a point $\mathbf{w} \in \mathcal{W}$ is defined as

$$\nabla_{\text{GN}}^2(\ell \circ f)(\mathbf{w}) := \partial f(\mathbf{w})^* \nabla^2 \ell(f(\mathbf{w})) \partial f(\mathbf{w}).$$

As studied in Section 17.2, the Gauss-Newton matrix is a key ingredient of the Gauss-Newton method. An advantage of the Gauss-Newton matrix is its positive semi-definiteness provided that ℓ is convex.

Proposition 9.1 (Positive semi-definiteness of the GN matrix). If ℓ is convex, then $\nabla_{\text{GN}}^2(\ell \circ f)(\mathbf{w})$ is positive semi-definite for all f .

This means that the approximation

$$L(\mathbf{w} + \mathbf{v}) \approx \langle \nabla L(\mathbf{w}), \mathbf{v} \rangle + \frac{1}{2} \langle \mathbf{v}, \nabla_{\text{GN}}^2 L(\mathbf{w}) \mathbf{v} \rangle$$

is a **convex quadratic** w.r.t. \mathbf{v} .

Using the chain rule, we find that the Hessian of $L = \ell \circ f$ decomposes into the sum of two terms (see also Proposition 9.7).

Proposition 9.2 (Approximation of the Hessian). For f differentiable and ℓ twice differentiable, we have

$$\begin{aligned}\nabla^2(\ell \circ f)(\mathbf{w}) &= \partial f(\mathbf{w})^* \nabla^2 \ell(f(\mathbf{w})) \partial f(\mathbf{w}) + \partial^2 f(\mathbf{w})^* [\nabla \ell(f(\mathbf{w}))] \\ &= \nabla_{\text{GN}}^2(\ell \circ f)(\mathbf{w}) + \sum_{j=1}^Z \nabla_j \ell(f(\mathbf{w})) \nabla^2 f_j(\mathbf{w}).\end{aligned}$$

If f is linear, then the Hessian and Gauss-Newton matrices coincide,

$$\nabla^2(\ell \circ f)(\mathbf{w}) = \nabla_{\text{GN}}^2(\ell \circ f)(\mathbf{w}).$$

The Gauss-Newton operator $\nabla_{\text{GN}}^2(\ell \circ f)$ can therefore be seen as an approximation of the Hessian $\nabla^2(\ell \circ f)$, with equality if f is linear.

9.2.2 Gauss-Newton chain rule

A chain rule for computing the Hessian of a composition of two functions is presented in Proposition 9.7, but the formula is relatively complicated, due to the cross-terms. In contrast, a Gauss-Newton chain rule is straightforward.

Proposition 9.3 (Gauss-Newton chain rule).

$$\nabla_{\text{GN}}^2(\ell \circ f \circ g)(\mathbf{w}) = \partial g(\mathbf{w})^* \nabla_{\text{GN}}^2(\ell \circ f)(g(\mathbf{w})) \partial g(\mathbf{w}).$$

9.2.3 Gauss-Newton vector product

As for the Hessian, we rarely need to materialize the full Gauss-Newton matrix in memory. Indeed, we can define the Gauss-Newton vector product (GNVP), a linear map for a direction $\mathbf{v} \in \mathcal{W}$, as

$$\nabla_{\text{GN}}^2(\ell \circ f)(\mathbf{w})[\mathbf{v}] := \partial f(\mathbf{w})^* \nabla^2 \ell(f(\mathbf{w})) \partial f(\mathbf{w}) \mathbf{v}, \quad (9.1)$$

where $\nabla^2 \ell(\boldsymbol{\theta}) \mathbf{u}$ is the HVP of ℓ , a linear map from \mathcal{M} to \mathcal{M} . The GNVP can be computed using the JVP of f , the HVP of ℓ and the VJP of f . Instantiating the VJP requires 1 forward pass through f , from which we get both the value $f(\mathbf{w})$ and the adjoint linear map $\mathbf{u} \mapsto (\partial f(\mathbf{w})^* \mathbf{u})$. Evaluating the VJP requires 1 backward pass through

f . Evaluating the JVP requires 1 forward pass through f . In total, evaluating $\mathbf{v} \mapsto \nabla_{\text{GN}}^2(\ell \circ f)(\mathbf{w})\mathbf{v}$ therefore requires 2 forward passes and 1 backward pass through f .

9.2.4 Gauss-Newton matrix factorization

In this section, we assume $\mathcal{W} \subseteq \mathbb{R}^P$ and $\mathcal{M} \subseteq \mathbb{R}^M$. When ℓ is convex, we know that the Gauss-Newton matrix is positive semi-definite and therefore it can be factorized into $\nabla_{\text{GN}}^2(\ell \circ f)(\mathbf{w}) = VV^\top$ for some $V \in \mathbb{R}^{P \times R}$, where $R \leq \min\{P, M\}$ is the rank of the matrix. Such a decomposition can actually be computed easily from a factorization of the Hessian of ℓ . For instance, suppose we know the eigendecomposition of the Hessian of ℓ , $\nabla^2 \ell(f(\mathbf{w})) = \sum_{j=1}^M \lambda_j \mathbf{u}_j \mathbf{u}_j^\top$, where the \mathbf{u}_j are the eigenvectors and the $\lambda_j \geq 0$ are the eigenvalues (which we know are non-negative due to positive semidefiniteness). Then, the Gauss-Newton matrix can be decomposed as

$$\begin{aligned} \nabla_{\text{GN}}^2(\ell \circ f) &= \sum_{j=1}^M \lambda_j \partial f(\mathbf{w})^* \mathbf{u}_j \mathbf{u}_j^\top \partial f(\mathbf{w})^* \\ &= \sum_{j=1}^M \left(\sqrt{\lambda_j} \partial f(\mathbf{w})^* \mathbf{u}_j \right) \left(\sqrt{\lambda_j} \partial f(\mathbf{w})^* \mathbf{u}_j \right)^\top \\ &= \sum_{j=1}^M \mathbf{v}_j \mathbf{v}_j^\top \quad \text{where } \mathbf{v}_j := \sqrt{\lambda_j} \partial f(\mathbf{w})^* \mathbf{u}_j. \end{aligned}$$

Stacking the vectors \mathbf{v}_j into a matrix $V = (\mathbf{v}_1, \dots, \mathbf{v}_M)$, we recover the factorization $\nabla_{\text{GN}}^2(\ell \circ f)(\mathbf{w}) = VV^\top$. To form this decomposition, we need to perform the eigendecomposition of $\nabla^2 \ell(f(\mathbf{w})) \in \mathbb{R}^{M \times M}$, which takes $O(M^3)$ time. We also need M calls to the VJP of f at \mathbf{w} . Compared to the direct implementation in Eq. (9.1), the factorization, once computed, allows us to compute the Gauss-Newton vector product (GNVP) as $\nabla_{\text{GN}}^2(\ell \circ f)(\mathbf{w})[\mathbf{v}] = VV^\top \mathbf{v}$. The factorization only requires $P \times M$ memory, while the direct implementation in Eq. (9.1) requires us to maintain the intermediate computations of f . The computation-memory trade-offs therefore depend on the function considered.

9.2.5 Stochastic setting

Suppose the objective function is of the form

$$L(\mathbf{w}; \mathbf{x}, \mathbf{y}) := \ell(f(\mathbf{w}; \mathbf{x}); \mathbf{y}).$$

With some slight abuse of notation, we then have that the Gauss-Newton matrix associated with a pair (\mathbf{x}, \mathbf{y}) is

$$\nabla_{\text{GN}}^2 L(\mathbf{w}; \mathbf{x}, \mathbf{y}) := \partial f(\mathbf{w}; \mathbf{x})^* \nabla^2 \ell(\boldsymbol{\theta}; \mathbf{y}) \partial f(\mathbf{w}; \mathbf{x}).$$

Given a distribution ρ over (\mathbf{x}, \mathbf{y}) pairs, the Gauss-Newton matrix associated with the averaged loss

$$L(\mathbf{w}) := \mathbb{E}_{X, Y \sim \rho} [L(\mathbf{w}; X, Y)]$$

is then

$$\nabla_{\text{GN}}^2 L(\mathbf{w}) = \mathbb{E}_{X, Y \sim \rho} \left[\nabla_{\text{GN}}^2 L(\mathbf{w}; X, Y) \right].$$

9.3 Fisher information matrix

9.3.1 Definition using the score function

The Fisher information is a way to measure the amount of information in a random variable S .

Definition 9.2 (Fisher information matrix). The **Fisher information matrix**, or Fisher for short, associated with the negative log-likelihood $L(\mathbf{w}; S) = -\log q_{\mathbf{w}}(S)$ of a probability distribution $q_{\mathbf{w}}$ with parameters \mathbf{w} is the covariance of the gradients of L at \mathbf{w} for S distributed according to $q_{\mathbf{w}}$,

$$\begin{aligned} \nabla_{\text{F}}^2 L(\mathbf{w}) &:= \mathbb{E}_{S \sim q_{\mathbf{w}}} [\nabla L(\mathbf{w}; S) \otimes \nabla L(\mathbf{w}; S)] \\ &= \mathbb{E}_{S \sim q_{\mathbf{w}}} [\nabla_{\mathbf{w}} \log q_{\mathbf{w}}(S) \otimes \nabla_{\mathbf{w}} \log q_{\mathbf{w}}(S)]. \end{aligned}$$

The gradient $\nabla_{\mathbf{w}} \log q_{\mathbf{w}}(S)$ is known as the **score function**.

As studied in Section 17.3, the Fisher information matrix is a key ingredient of the natural gradient descent method.

9.3.2 Link with the Hessian

Provided that the probability distribution is twice differentiable w.r.t. \mathbf{w} with integrable second derivatives, the Fisher information matrix can also be expressed as the Hessian of the negative log-likelihood (Amari, 1998; Martens, 2020).

Proposition 9.4 (Connection with the Hessian). The Fisher information matrix of the negative log-likelihood $L(\mathbf{w}; S) = -\log q_{\mathbf{w}}(S)$ satisfies

$$\nabla_{\mathbf{F}}^2 L(\mathbf{w}) = \mathbb{E}_{S \sim q_{\mathbf{w}}} [\nabla^2 L(\mathbf{w}; S)] = \mathbb{E}_{S \sim q_{\mathbf{w}}} [-\nabla_{\mathbf{w}}^2 \log q_{\mathbf{w}}(S)].$$

Remark 9.1 (Empirical Fisher). We emphasize that in the above definitions, S is sampled from the model distribution $q_{\mathbf{w}}$, not from the data distribution ρ . That is, we have

$$\begin{aligned} \nabla_{\mathbf{F}}^2 L(\mathbf{w}) &= \mathbb{E}_{S \sim q_{\mathbf{w}}} [\nabla_{\mathbf{w}} \log q_{\mathbf{w}}(S) \nabla_{\mathbf{w}} \log q_{\mathbf{w}}(S)^{\top}] \\ &\neq \mathbb{E}_{S \sim \rho} [\nabla_{\mathbf{w}} \log q_{\mathbf{w}}(S) \nabla_{\mathbf{w}} \log q_{\mathbf{w}}(S)^{\top}] \end{aligned}$$

The latter is sometimes called ambiguously the “empirical” Fisher, though this name has generated confusion (Kunstner *et al.*, 2019).

9.3.3 Equivalence with the Gauss-Newton matrix

So far, we discussed the Fisher information for a generic random variable $S \sim q_{\mathbf{w}}$. We now discuss the supervised probabilistic learning setting where $S = (X, Y)$ and where, using the product rule of probability, we define the PDF $q_{\mathbf{w}}(X, Y) := \rho_X(X)p_{\boldsymbol{\theta}}(Y)$, with the shorthand $\boldsymbol{\theta} := f(\mathbf{w}; X)$.

Proposition 9.5 (Fisher matrix in supervised setting). Suppose $(X, Y) \sim q_{\mathbf{w}}$ where the PDF of $q_{\mathbf{w}}$ is $q_{\mathbf{w}}(X, Y) := \rho_X(X)p_{\boldsymbol{\theta}}(Y)$. In that case, the Fisher information matrix of the negative log-

likelihood $L(\mathbf{w}; \mathbf{x}, \mathbf{y}) = -\log q_{\mathbf{w}}(\mathbf{x}, \mathbf{y})$ decomposes as,

$$\begin{aligned}\nabla_{\mathbf{F}}^2 L(\mathbf{w}) &= \mathbb{E}_{(X,Y) \sim q_{\mathbf{w}}} [\nabla_{\mathbf{w}} \log q_{\mathbf{w}}(X, Y) \otimes \nabla_{\mathbf{w}} \log q_{\mathbf{w}}(X, Y)] \\ &= \mathbb{E}_{X \sim \rho_X} [\mathbb{E}_{Y \sim p_{\theta}} [\nabla_{\mathbf{w}} \log p_{\theta}(Y) \otimes \nabla_{\mathbf{w}} \log p_{\theta}(Y)]] \\ &= \mathbb{E}_{X \sim \rho_X} [\partial f(\mathbf{w}; X)^* \nabla_{\mathbf{F}}^2 \ell(\theta) \partial f(\mathbf{w}; X)],\end{aligned}$$

where we defined the shorthand $\theta := f(\mathbf{w}; X)$ and where we defined the negative log-likelihood loss $\ell(\theta; Y) := -\log p_{\theta}(Y)$.

When p_{θ} is an exponential family distribution, we can show that the Fisher information matrix and the Gauss-Newton matrix are equivalent.

Proposition 9.6 (Equivalence between Fisher and Gauss-Newton). If p_{θ} is an exponential family distribution, then

$$\begin{aligned}\nabla_{\mathbf{F}}^2 L(\mathbf{w}) &= \mathbb{E}_{X \sim \rho_X} \mathbb{E}_{Y \sim p_{\theta}} [\nabla L(\mathbf{w}; X, Y) \otimes \nabla L(\mathbf{w}; X, Y)] \\ &= \mathbb{E}_{X \sim \rho_X} \mathbb{E}_{Y \sim p_{\theta}} [\partial f(\mathbf{w}; X)^* \nabla \ell(\theta; Y) \otimes \nabla \ell(\theta; Y) \partial f(\mathbf{w}; X)] \\ &= \mathbb{E}_{X \sim \rho_X} \mathbb{E}_{Y \sim p_{\theta}} [\partial f(\mathbf{w}; X)^* \nabla^2 \ell(\theta; Y) \partial f(\mathbf{w}; X)] \\ &= \mathbb{E}_{X, Y \sim \rho} [\nabla_{\text{GN}}^2 L(\mathbf{w}; X, Y)],\end{aligned}$$

where $\rho_X(\mathbf{x}) := \int \rho(\mathbf{x}, \mathbf{y}) d\mathbf{y}$.

Proof. From Proposition 3.3, if p_{θ} is an exponential family distribution, $\nabla^2 \ell(\theta, \mathbf{y})$ is actually independent of \mathbf{y} . Using Bartlett's second identity Eq. (12.3), we then obtain

$$\begin{aligned}\nabla^2 \ell(\theta; \cdot) &= \mathbb{E}_{Y \sim p_{\theta}} [\nabla^2 \ell(\theta; Y)] \\ &= \mathbb{E}_{Y \sim p_{\theta}} [\nabla^2 \ell(\theta; Y)] \\ &= \mathbb{E}_{Y \sim p_{\theta}} [-\nabla_{\theta}^2 \log p_{\theta}(Y)] \\ &= \mathbb{E}_{Y \sim p_{\theta}} [\nabla_{\theta} \log p_{\theta}(Y) \otimes \nabla_{\theta} \log p_{\theta}(Y)] \\ &= \mathbb{E}_{Y \sim p_{\theta}} [\nabla \ell(\theta; Y) \otimes \nabla \ell(\theta; Y)],\end{aligned}$$

where we used \cdot to indicate that the results holds for all \mathbf{y} . Plugging the result back in the Fisher information matrix concludes the proof. \square

9.4 Inverse-Hessian vector product

9.4.1 Definition as a linear map

We saw in Section 17.1 that Newton's method uses iterations as

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \nabla^2 L(\mathbf{w}^t)^{-1} \nabla L(\mathbf{w}^t).$$

The inverse is well-defined if for example L is strictly convex. Otherwise, we saw that some additional regularization can be added. Newton's method therefore requires to access **inverse-Hessian vector products** (IHVPs), as defined below.

Definition 9.3 (Inverse-Hessian vector product). For a twice differentiable function $L : \mathbb{R}^P \rightarrow \mathbb{R}$, we define the **inverse-Hessian Vector Product** (IHVP) of L at $\mathbf{w} \in \mathbb{R}^P$ as the linear map

$$\mathbf{u} \mapsto \nabla^2 L(\mathbf{w})^{-1} \mathbf{u},$$

provided that it exists. In other words, it is the linear map which to \mathbf{u} associates \mathbf{v} such that $\nabla^2 L(\mathbf{w})\mathbf{v} = \mathbf{u}$.

9.4.2 Implementation with matrix-free linear solvers

Numerous direct methods exist to compute the inverse of a matrix, such as the Cholesky decomposition, QR decomposition and Gaussian elimination. However, these algorithms require accessing elementary entries of the matrix, while an autodiff framework gives access to the Hessian through HVPs. Fortunately, there exists so-called **matrix-free** algorithms, that can solve a linear system of equations

$$H[\mathbf{v}] = \mathbf{u}$$

by only accessing the linear map $\mathbf{v} \mapsto H[\mathbf{v}]$ for any \mathbf{v} . Among such algorithms, we have the **conjugate gradient** (CG) method, that applies for H positive-definite, i.e., such that $\langle \mathbf{v}, H[\mathbf{v}] \rangle > 0$ for all $\mathbf{v} \neq 0$, or the **generalized minimal residual** (GMRES) method, that applies for any invertible H . A longer list of solvers can be found in public software such as SciPy (Virtanen *et al.*, 2020). The IHVP of a strictly convex

function (ensuring that the Hessian is positive definite) can therefore be computed by instantiating CG on the HVP,

$$\nabla^2 L(\mathbf{w})^{-1} \mathbf{u} \approx \text{CG}(\nabla^2 L(\mathbf{w})[\cdot], \mathbf{u}).$$

Positive-definiteness of the Hessian is indeed guaranteed for strictly convex functions for example, while for generic non-convex functions, such property may be verified around a minimizer but not in general. The conjugate gradient method is recalled in Algorithm 9.1 in its simplest form. In theory, the exact solution of the linear system is found after at most $T = P$ iterations of CG, though in practice numerical errors may prevent from getting an exact solution.

Algorithm 9.1 Conjugate gradient method

Inputs: linear map $H[\cdot] : \mathbb{R}^P \rightarrow \mathbb{R}^P$, target $\mathbf{u} \in \mathbb{R}^P$, initialization \mathbf{v}_0 (default $\mathbf{0}$), number of iterations T (default P), target accuracy ε (default machine precision)

- 1: $\mathbf{r}_0 = \mathbf{u} - H[\mathbf{v}_0]$
- 2: $\mathbf{p}_0 = \mathbf{r}_0$
- 3: **for** $t = 0, \dots, T$ **do**
- 4: $\alpha_t = \frac{\langle \mathbf{r}_t, \mathbf{r}_t \rangle}{\langle \mathbf{p}_t, H[\mathbf{p}_t] \rangle}$
- 5: $\mathbf{v}_{t+1} = \mathbf{v}_t + \alpha_t \mathbf{p}_t$
- 6: $\mathbf{r}_{t+1} = \mathbf{r}_t - \alpha_t H[\mathbf{p}_t]$
- 7: **if** $\langle \mathbf{r}_{t+1}, \mathbf{r}_{t+1} \rangle \leq \varepsilon$ **then break**
- 8: $\beta_t = \frac{\langle \mathbf{r}_{t+1}, \mathbf{r}_{t+1} \rangle}{\langle \mathbf{r}_t, \mathbf{r}_t \rangle}$
- 9: $\mathbf{p}_{t+1} = \mathbf{r}_{t+1} + \beta_t \mathbf{p}_t$

Output: \mathbf{v}_T , such that $H[\mathbf{v}_T] \approx \mathbf{u}$

9.4.3 Complexity

For a given matrix $H \in \mathbb{R}^{P \times P}$, solving $H\mathbf{v} = \mathbf{u}$ can be done with decomposition methods (LU, QR, Cholesky) in $O(P^3)$ time. For matrix-free methods such as CG or GMRES, the cost per iteration is $O(P^2)$. Since they theoretically solve the linear system in $O(P)$ iterations, the cost to obtain an exact solution is theoretically the same, $O(P^3)$.

However, CG or GMRES differ from decomposition methods in that they are iterative methods, meaning that, at each iteration, they get closer to a solution. Unlike decomposition methods, this means that we can stop them before an exact solution is found. In practice, the number of iterations required to find a good approximate solution depends on the matrix. Well conditioned matrices require only few iterations. Badly conditioned matrices lead to some numerical instabilities for CG, so that more than P iterations may be needed to get a good solution. In contrast, decomposition methods proceed in two steps: first they build a decomposition of H at a cost of $O(P^3)$, and second they solve a linear system at a cost of $O(P^2)$, by leveraging the structure. LU and QR decompositions are known to be generally more stable and are therefore often preferred in practice, when we can access entries of H at no cost.

If we do not have access to the Hessian H , but only to its HVP, accessing entries of H comes at a prohibitive cost. Indeed, entries of H can still be recovered from HVPs, since $\mathbf{e}_i^\top H \mathbf{e}_j = H_{i,j}$, but accessing each row or column of H costs one HVP (matrix-vector product). To access the information necessary to use a decomposition method, we therefore need P calls to HVPs before being able to actually compute the solution. For the same number of calls, CG or GMRES will already have found an approximate solution. In addition, a CG method does not require to store any memory.

9.5 Second-order backpropagation

9.5.1 Second-order Jacobian chain rule

The essential ingredient to develop forward-mode and reverse-mode autodiff hinged upon the chain rule for composed functions, $h = g \circ f$. For second derivatives, a similar rule can be obtained. To do so, we slightly abuse notations and denote

$$\partial^2 h(\mathbf{w})^*[\mathbf{u}] := \nabla^2 \langle h, \mathbf{u} \rangle(\mathbf{w}) \in \mathbb{R}^{P \times P},$$

where $h : \mathbb{R}^P \rightarrow \mathbb{R}^Q$, $\mathbf{w} \in \mathbb{R}^P$, $\mathbf{u} \in \mathbb{R}^Q$, and where we recall the shorthand notation $\langle \mathbf{u}, h \rangle(\mathbf{w}) := \langle \mathbf{u}, h(\mathbf{w}) \rangle$. Moreover, we view the above quantity as a linear map. Strictly speaking, the superscript $*$ is

not a linear adjoint anymore, since $\mathbf{v}_1, \mathbf{v}_2 \mapsto \partial^2 h(\mathbf{w})[\mathbf{v}_1, \mathbf{v}_2]$ is no longer linear but bilinear. However, this superscript plays the same role as the VJP, since it takes an output vector and returns the input derivatives that correspond to infinitesimal variations along that output vector.

Proposition 9.7 (Hessian chain-rule). For two twice differentiable functions $f: \mathbb{R}^P \rightarrow \mathbb{R}^M$ and $g: \mathbb{R}^M \rightarrow \mathbb{R}^Q$, the second directional derivative of the composition $g \circ f$ is a bilinear map from $\mathbb{R}^P \times \mathbb{R}^P$ to \mathbb{R}^Q along input directions $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^P$ of the form

$$\begin{aligned} \partial^2(g \circ f)(\mathbf{w})[\mathbf{v}_1, \mathbf{v}_2] &= \partial g(f(\mathbf{w}))[\partial^2 f(\mathbf{w})[\mathbf{v}_1, \mathbf{v}_2]] \\ &\quad + \partial^2 g(f(\mathbf{w}))[\partial f(\mathbf{w})[\mathbf{v}_1], \partial f(\mathbf{w})[\mathbf{v}_2]]. \end{aligned}$$

The Hessian of the composition $g \circ f$ along an output direction $\mathbf{u} \in \mathbb{R}^Q$ is, seen as a linear map,

$$\begin{aligned} \partial^2(g \circ f)(\mathbf{w})^*[\mathbf{u}] &= \partial^2 f(\mathbf{w})^*[\partial g(f(\mathbf{w}))^* \mathbf{u}] \\ &\quad + \partial f(\mathbf{w})^* \partial^2 g(f(\mathbf{w}))^*[\mathbf{u}] \partial f(\mathbf{w}). \end{aligned} \tag{9.2}$$

For the composition of $f: \mathbb{R}^P \rightarrow \mathbb{R}^M$ with a scalar-valued function $\ell: \mathbb{R}^M \rightarrow \mathbb{R}$, we have in matrix form

$$\begin{aligned} \nabla^2(\ell \circ f)(\mathbf{w}) &= \sum_{j=1}^M (\nabla \ell(f(\mathbf{w})))_j \nabla^2 f_j(\mathbf{w}) \\ &\quad + \partial f(\mathbf{w})^\top \nabla^2 \ell(f(\mathbf{w})) \partial f(\mathbf{w}). \end{aligned}$$

Note that, while the Hessian is usually defined for scalar-valued functions $h: \mathbb{R}^P \rightarrow \mathbb{R}$, the above definition is for a generalized notion of Hessian that works for any function $h: \mathbb{R}^P \rightarrow \mathbb{R}^Q$.

The Hessian back-propagation rule in Eq. (9.2) reveals two terms. The first one $\partial^2 f(\mathbf{w})^*[\partial g(f(\mathbf{w}))^* \mathbf{u}]$ simply computes the Hessian of the intermediate function along the output direction normally back-propagated by a VJP. The second term $\partial f(\mathbf{w})^* \partial^2 g(f(\mathbf{w}))^*[\mathbf{u}] \partial f(\mathbf{w})$ shows how intermediate first-order variations influence second order derivatives of the output.

Example 9.1 (Composition with an elementwise nonlinear function).

Consider the element-wise application of a twice differentiable scalar-valued function $f(\mathbf{x}) = (f(x_i))_{i=1}^M$ followed by some twice differentiable function ℓ . Note that $\nabla^2 f_i(\mathbf{x}) = f''(x_i) \mathbf{e}_i \mathbf{e}_i^\top$. Hence, the Hessian of the composition reads

$$\begin{aligned} \nabla^2(\ell \circ f)(\mathbf{x}) &= \sum_{i=1}^M (\nabla \ell(f(\mathbf{x})))_i f''(x_i) \mathbf{e}_i \mathbf{e}_i^\top \\ &\quad + \mathbf{diag}(f'(\mathbf{x})) \nabla^2 \ell(f(\mathbf{x})) \mathbf{diag}(f'(\mathbf{x})) \\ &= \mathbf{diag}(\nabla \ell(f(\mathbf{w}))) \odot f''(\mathbf{x}) \\ &\quad + \nabla^2 \ell(f(\mathbf{x})) \odot (f'(\mathbf{x}) f'(\mathbf{x})^\top), \end{aligned}$$

where $f'(\mathbf{x}) := (f'(x_i))_{i=1}^M$ and $f''(\mathbf{x}) := (f''(x_i))_{i=1}^M$.

Example 9.2 (Hessian of the composition with a linear function). Consider

a linear function $f(\mathbf{W}) = \mathbf{W}\mathbf{x}$, for $\mathbf{W} \in \mathbb{R}^{M \times D}$, composed with some twice differentiable function $\ell : \mathbb{R}^M \rightarrow \mathbb{R}$. From Proposition 9.7, we get, in terms of linear maps,

$$\nabla^2(\ell \circ f)(\mathbf{W}) = \partial f(\mathbf{W})^* \nabla^2 \ell(f(\mathbf{W})) \partial f(\mathbf{W}).$$

As already noted in Section 2.3, we have that $\partial f(\mathbf{W})[\mathbf{V}] = \mathbf{V}\mathbf{x}$ and $\partial f(\mathbf{W})^*[\mathbf{u}] = \mathbf{u}\mathbf{x}^\top$. Hence, the Hessian seen as a linear map reads

$$\nabla^2(\ell \circ f)(\mathbf{W})[\mathbf{V}] = \partial f(\mathbf{W})^*[\nabla^2 \ell(f(\mathbf{W}))[\partial f(\mathbf{W})[\mathbf{V}]]] = \mathbf{H}\mathbf{V}\mathbf{x}\mathbf{x}^\top,$$

where $\mathbf{H} := \nabla^2 \ell(f(\mathbf{W}))$.

9.5.2 Computation chains

For a simple computation chain $f = f_K \circ \dots \circ f_1$ as in Section 8.1, the formula derived in Proposition 9.7 suffices to develop an algorithm that backpropagates the Hessian, as shown in Algorithm 9.2. Compared to Algorithm 8.2, we simply backpropagate both the vectors \mathbf{r}_k and the matrices \mathbf{R}_k using intermediate first and second derivatives.

Algorithm 9.2 Hessian backprop for computation chains**Functions:** $f := f_K \circ \dots \circ f_1$,**Inputs:** input \mathbf{x} , output direction \mathbf{u}

- 1: Initialize and store $\mathbf{s}_0 := \mathbf{x}$ ▷ Forward pass
- 2: **for** $k := 1, \dots, K$ **do**
- 3: Compute and store $\mathbf{s}_k := f_k(\mathbf{s}_{k-1})$
- 4: Initialize $\mathbf{r}_K := \nabla \ell(\mathbf{s}_K)$, $\mathbf{R}_K := \nabla^2 \ell(\mathbf{s}_K)$ ▷ Backward pass
- 5: **for** $k := K, \dots, 1$ **do**
- 6: Compute $\mathbf{r}_{k-1} := \partial f_k(\mathbf{s}_{k-1})^* [\mathbf{r}_k]$
- 7: Compute $\mathbf{R}_{k-1} := \partial^2 f_k(\mathbf{s}_{k-1})^* [\mathbf{r}_k] + \partial f_k(\mathbf{s}_{k-1})^* \mathbf{R}_k \partial f_k(\mathbf{s}_{k-1})$
- 8: Release \mathbf{s}_{k-1} from memory

Outputs: $\ell(f(\mathbf{x})) = \ell(\mathbf{s}_K)$, $\nabla(\ell \circ f)(\mathbf{x}) = \mathbf{r}_0$, $\nabla^2(\ell \circ f)(\mathbf{x}) = \mathbf{R}_0$ **9.5.3 Fan-in and fan-out**

For generic computation graphs (see Section 8.3), we saw that multi-input functions (fan-in) were crucial. For Hessian backpropagation in computation graphs, we therefore need to develop a similar formula.

Proposition 9.8 (Hessian chain-rule for fan-in). Consider $n+1$ twice differentiable functions f_1, \dots, f_n and g with $f_i : \mathbb{R}^P \rightarrow \mathbb{R}^{M_i}$ and $g : \mathbb{R}^{M_1} \times \dots \times \mathbb{R}^{M_n} \rightarrow \mathbb{R}^Q$. The Hessian of $g \circ f$ for $f(\mathbf{w}) = (f_1(\mathbf{w}), \dots, f_n(\mathbf{w}))$ along an output direction $\mathbf{u} \in \mathbb{R}^Q$ is given by

$$\begin{aligned} \partial^2(g \circ f)(\mathbf{w})^* [\mathbf{u}] &= \sum_{i=1}^n \partial^2 f_i(\mathbf{w})^* [\partial_i g(f(\mathbf{w}))^* [\mathbf{u}]] \\ &\quad + \sum_{i,j=1}^n \partial f_i(\mathbf{w})^* \partial_{i,j}^2 g(f(\mathbf{w}))^* [\mathbf{u}] \partial f_j(\mathbf{w}). \end{aligned}$$

The gradient backpropagation expression for fan-in is simple because the functions f_i are not linked by any path. In contrast, the Hessian backpropagation involves cross-product terms

$\partial f_i(\mathbf{w})^* \partial_{i,j}^2 g(f(\mathbf{w}))^* [\mathbf{u}] \partial f_j(\mathbf{w})$ for $i \neq j$. The nodes associated to the f_i computations cannot be treated independently anymore.

On the other hand, developing a backpropagation rule for fan-out does not pose any issue, since each output function can be treated

independently.

Proposition 9.9 (Hessian chain-rule for fan-out). Consider $n+1$ twice differentiable functions g_1, \dots, g_n and f with $g_i : \mathbb{R}^M \rightarrow \mathbb{R}^{Q_i}$ and $f : \mathbb{R}^P \rightarrow \mathbb{R}^M$. The Hessian of $g \circ f$ for $g(\mathbf{w}) = (g_1(\mathbf{w}), \dots, g_n(\mathbf{w}))$ along a direction $\mathbf{u} = (\mathbf{u}_1, \dots, \mathbf{u}_n) \in \mathbb{R}^{Q_1} \times \dots \times \mathbb{R}^{Q_n}$ is given by

$$\begin{aligned} \partial^2(g \circ f)(\mathbf{w})^*[\mathbf{u}] &= \sum_{i=1}^n \partial^2 f(\mathbf{w})^*[\partial g_i(f(\mathbf{w}))^*[\mathbf{u}_i]] \\ &\quad + \sum_{i=1}^n \partial f(\mathbf{w})^* \partial^2 g_i(f(\mathbf{w}))^*[\mathbf{u}] \partial f(\mathbf{w}). \end{aligned}$$

9.6 Block diagonal approximations

Rather than computing the whole Hessian or Gauss-Newton matrices, we can consider computing block-diagonal or diagonal approximations, which are easier to invert. The approximation rules we present in this section build upon the Hessian chain rule studied in Section 9.5.

9.6.1 Feedforward networks

Recall the definition of a feedforward network:

$$\begin{aligned} \mathbf{s}_0 &:= \mathbf{x} \\ \mathbf{s}_k &:= f_k(\mathbf{s}_{k-1}, \mathbf{w}_k) \quad \forall k \in \{1, \dots, K\} \\ f(\mathbf{x}, \mathbf{w}) &:= \mathbf{s}_K, \end{aligned}$$

where $\mathbf{w} := (\mathbf{w}_1, \dots, \mathbf{w}_K)$. Rather than computing the entire Hessian of $\ell \circ f$ w.r.t. \mathbf{w} , we can compute the Hessians w.r.t. each set of parameters \mathbf{w}_k . For the case of computation chains, the Hessian backpropagation recursion we used in Algorithm 9.2 was

$$\mathbf{R}_{k-1} := \partial^2 f_k(\mathbf{s}_{k-1})^*[\mathbf{r}_k] + \partial f_k(\mathbf{s}_{k-1})^* \mathbf{R}_k \partial f_k(\mathbf{s}_{k-1}).$$

Extending this recursion to the feedforward network case, we obtain, starting from $\mathbf{r}_K := \nabla \ell(\mathbf{s}_K)$ and $\mathbf{R}_K := \nabla^2 \ell(\mathbf{s}_K)$,

$$\begin{aligned} \mathbf{r}_{k-1} &:= \partial f_k(\mathbf{s}_{k-1}, \mathbf{w}_k)^* [\mathbf{r}_k] \\ \begin{pmatrix} \mathbf{R}_{k-1} & \sim \\ \sim & \mathbf{H}_k \end{pmatrix} &:= \partial^2 f_k(\mathbf{s}_{k-1}, \mathbf{w}_k)^* [\mathbf{r}_k] \\ &\quad + \partial f_k(\mathbf{s}_{k-1}, \mathbf{w}_k)^* \mathbf{R}_k \partial f_k(\mathbf{s}_{k-1}, \mathbf{w}_k), \end{aligned}$$

where we used \sim to indicate that these blocks are not used. The Hessians w.r.t each set of parameters are then

$$\begin{aligned} \mathbf{R}_0 &= \nabla_{\mathbf{x}\mathbf{x}}^2 (\ell \circ f)(\mathbf{x}, \mathbf{w}) \\ \mathbf{H}_1 &= \nabla_{\mathbf{w}_1 \mathbf{w}_1}^2 (\ell \circ f)(\mathbf{x}, \mathbf{w}) \\ &\vdots \\ \mathbf{H}_K &= \nabla_{\mathbf{w}_K \mathbf{w}_K}^2 (\ell \circ f)(\mathbf{x}, \mathbf{w}). \end{aligned}$$

The validity of this result stems from the fact that we can view the Hessian w.r.t. \mathbf{w}_k as computing the Hessian w.r.t. \mathbf{w}_k of

$$\tilde{f}_K \circ \dots \circ \tilde{f}_{k+1} \circ f_k(\mathbf{s}_{k-1}, \mathbf{w}_k)$$

where $\tilde{f}_i := f_i(\cdot, \mathbf{w}_k)$, for $i \in \{k+1, \dots, K\}$. As the computations of the block-wise Hessians share most of the computations, they can be evaluated in a single backward pass just as the gradients.

Example 9.3 (Block-wise computation of the Gauss-Newton matrix). Our blockwise backpropagation scheme can readily be adapted for the Gauss-Newton matrix as

$$\begin{pmatrix} \mathbf{R}_{k-1} & \sim \\ \sim & \mathbf{G}_k \end{pmatrix} := \partial f_k(\mathbf{s}_{k-1}, \mathbf{w}_k)^* \mathbf{R}_k \partial f_k(\mathbf{s}_{k-1}, \mathbf{w}_k),$$

starting from $\mathbf{R}_K := \nabla^2 \ell(\mathbf{s}_K)$. The outputs $\mathbf{R}_0, \mathbf{G}_0, \dots, \mathbf{G}_K$ give a block-wise approximation of the Gauss-Newton matrix.

Now, consider a simple multilayer perceptron such that

$$f_k(\mathbf{s}_{k-1}, \mathbf{w}_k) := \mathbf{a}(\mathbf{W}_k \mathbf{s}_{k-1}) \quad \text{with} \quad \mathbf{w}_k := \text{vec}(\mathbf{W}_k)$$

Using Example 9.2 and Example 9.1 adapted to the Gauss-Newton

matrix, we can compute the block-wise decomposition of the Gauss-Newton matrix as, for $k = K, \dots, 1$,

$$\begin{aligned}\mathbf{R}_{k-1} &:= \mathbf{W}_k^\top \mathbf{J}_k \mathbf{W}_k \\ \mathbf{J}_k &:= \mathbf{R}_k \odot (\mathbf{a}'(\mathbf{W}_k \mathbf{s}_{k-1}) \mathbf{a}'(\mathbf{W}_k \mathbf{s}_{k-1})^\top) \\ \mathbf{G}_k &:= \mathbf{J}_k \otimes \mathbf{s}_{k-1} \mathbf{s}_{k-1}^\top\end{aligned}$$

starting from $\mathbf{R}_K := \nabla^2 \ell(\mathbf{s}_K)$. The outputs $\mathbf{G}_1, \dots, \mathbf{G}_K$ correspond to the block-wise elements of the Gauss-Newton matrix of f for the vectorized weights $\mathbf{w}_1, \dots, \mathbf{w}_K$. Similar computations were done in KFRA (Botev *et al.*, 2017) and BackPack (Dangel *et al.*, 2019).

9.6.2 Computation graphs

For generic computation graphs, consider a function $f(\mathbf{x}, \mathbf{w})$ defined by, denoting $i_1, \dots, i_{p_k} := \text{pa}(k)$,

$$\mathbf{s}_k := f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}}) \quad \forall k \in \{1, \dots, K\}$$

such that $f(\mathbf{x}, \mathbf{w}) = \mathbf{s}_K$, and k is following a topological ordering of the graph (see Section 8.3). We can consider the following backpropagation scheme, for $k = K, \dots, 1$ and $j \in \text{pa}(k)$

$$\mathbf{r}_{i_j} \leftarrow \mathbf{r}_{i_j} + \partial_j f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}})^* [\mathbf{r}_k] \quad (9.3)$$

$$\begin{aligned}R_{i_j} &\leftarrow R_{i_j} + \partial_{jj}^2 f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}})^* [\mathbf{r}_k] \\ &\quad + \partial_j f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}})^* R_k \partial_j f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}}),\end{aligned} \quad (9.4)$$

starting from $R_K := \nabla^2 \ell(\mathbf{s}_K)$ and $\mathbf{r}_K := \nabla \ell(\mathbf{s}_K)$. Recall that for multiple inputs, the chain-rule presented in Proposition 9.8 involves the cross-derivatives. For this reason the back-propagation scheme in Eq. (9.3) only computes an approximation. For example, one can verify that using Eq. (9.3) to compute the Hessian of $\ell(f_1(\mathbf{w}), f_2(\mathbf{w}))$ does not provide an exact expression for the Hessian of f . This scheme is easy to implement and may provide a relevant proxy for the Hessian.

9.7 Diagonal approximations

Similarly to the idea of designing a backpropagation scheme that approximates blocks of the Hessian, we can design a backpropagation

scheme that approximates the diagonal of the Hessian. The approach was originally proposed by Becker and Le Cun (1988) for feedforward networks, but our exposition, new to our knowledge, has the benefit that it naturally extends to computational graphs, as we shall see.

9.7.1 Computation chains

The idea stems from modifying the Hessian backpropagation rule in Proposition 9.7 to only keep the diagonal of the Hessian. Formally, given a matrix $\mathbf{M} \in \mathbb{R}^{D \times D}$, we denote by $\mathbf{diag}(\mathbf{M}) = (\mathbf{M}_{ii})_{i=1}^D \in \mathbb{R}^D$ the vector of diagonal entries of \mathbf{M} , and for a vector $\mathbf{m} \in \mathbb{R}^D$, we denote $\mathbf{Diag}(\mathbf{m}) = \sum_{i=1}^D m_i \mathbf{e}_i \mathbf{e}_i^\top$ the diagonal matrix with entries m_i . For the backpropagation of the Hessian of $\ell \circ f_K \circ \dots \circ f_1$, we see from Algorithm 9.2 that $\mathbf{diag}(\mathbf{H}_{k-1})$ can be expressed in terms of \mathbf{H}_k as

$$\begin{aligned} \mathbf{diag}(\mathbf{H}_{k-1}) &= \mathbf{diag}(\partial^2 f_k(\mathbf{s}_{k-1})^* \mathbf{r}_k) \\ &\quad + \mathbf{diag}(\partial f_k(\mathbf{s}_{k-1})^* \mathbf{H}_k \partial f_k(\mathbf{s}_{k-1})). \end{aligned}$$

Unfortunately, that recursion needs access to the whole Hessian \mathbf{H}_k , and would therefore be too expensive. A natural idea is to modify the recursion to approximate $\mathbf{diag}(\mathbf{H}_k)$ by backpropagating **vectors**:

$$\begin{aligned} \mathbf{d}_{k-1} &:= \mathbf{diag}(\partial^2 f_k(\mathbf{s}_{k-1})^* \mathbf{r}_k) \\ &\quad + \mathbf{diag}(\partial f_k(\mathbf{s}_{k-1})^* \mathbf{Diag}(\mathbf{d}_k) \partial f_k(\mathbf{s}_{k-1})). \end{aligned}$$

The diagonal matrix $\mathbf{Diag}(\mathbf{d}_k)$ serves as a surrogate for \mathbf{H}_k . Each iteration of this recursion can be computed in linear time in the output dimension D_k since

$$d_{k-1,i} = \sum_{j=1}^{D_k} r_{k,j} \cdot \partial_{i,i}^2 f_{k,j}(\mathbf{s}_{k-1}) + \sum_{j=1}^{D_k} d_{k,j} (\partial_i f_{k,j}(\mathbf{s}_{k-1}))^2.$$

To initialize the recursion, we can set $\mathbf{d}_K := \mathbf{diag}(\nabla^2 \ell(\mathbf{s}_K))$. As an alternative, as proposed by Elsayed and Mahmood (2022), if \mathbf{H}_K has a simple form, we can use $\nabla^2 \ell(\mathbf{s}_K)$ instead of $\mathbf{Diag}(\mathbf{d}_K)$ at the first iteration. This is the case for instance if f_K is a cross-entropy loss. The recursion is repeated until we obtain the approximate diagonal Hessian $\mathbf{d}_0 \approx \mathbf{diag}(\nabla^2(\ell \circ f)(\mathbf{x}))$. The gradients \mathbf{r}_k , needed to compute \mathbf{d}_k , are computed along the way and the algorithm can therefore also return $\mathbf{r}_0 = \nabla(\ell \circ f)(\mathbf{x})$.

9.7.2 Computation graphs

Although this diagonal approximation was originally derived for feed-forward networks Becker and Le Cun (1988), it is straightforward to generalize it to computation graphs. Namely, for a function $f(\mathbf{x}, \mathbf{w})$ decomposed along a computation graph, we can backpropagate a diagonal approximation in reverse topological order as

$$\begin{aligned} \mathbf{r}_{i_j} &\leftarrow \mathbf{r}_{i_j} + \partial_j f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}})^* [\mathbf{r}_k] \\ \mathbf{d}_{i_j} &\leftarrow \mathbf{d}_{i_j} + \mathbf{diag}(\partial_{jj}^2 f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}})^* [\mathbf{r}_k]) \\ &\quad + \mathbf{diag}(\partial_j f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}})^* \mathbf{Diag}(\mathbf{d}_k) \partial_j f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}})), \end{aligned} \quad (9.5)$$

for $j \in \text{pa}(k)$, starting from $\mathbf{r}_K = \nabla \ell(\mathbf{s}_K)$ and $\mathbf{d}_K = \mathbf{diag}(\nabla^2 \ell(\mathbf{s}_K))$ or $\mathbf{Diag}(\mathbf{d}_K) = \nabla^2 \ell(\mathbf{s}_K)$. To implement such an algorithm, each elementary function in the computational graph needs to be augmented with an oracle that computes the Hessian diagonal approximation of the current function, given the previous ones. An example with MLPs is presented in Example 9.4.

Example 9.4 (Hessian diagonal approximation for MLPs). Consider a multilayer perceptron

$$\begin{aligned} \mathbf{s}_k &:= \mathbf{a}_k(\mathbf{W}_k \mathbf{s}_{k-1}) \quad \forall k \in \{1, \dots, K-1\} \\ f(\mathbf{w}, \mathbf{x}) &:= \mathbf{s}_K \end{aligned}$$

starting from $\mathbf{s}_0 = \mathbf{x}$. Here \mathbf{a}_k is the element-wise activation function (potentially the identity) and \mathbf{w} encapsulates the weight matrices $\mathbf{W}_1, \dots, \mathbf{W}_K$. We consider the derivatives w.r.t. the flattened matrices, so that gradients and diagonal approximations w.r.t. these flattened quantities are vectors. The backpropagation scheme (9.5)

then reduces to, denoting $\mathbf{t}_k = \mathbf{W}_k \mathbf{s}_{k-1}$,

$$\begin{aligned}\mathbf{r}_{k-1} &:= \mathbf{W}_k^\top (\mathbf{a}'(\mathbf{t}_k) \odot \mathbf{r}_k) \\ \mathbf{g}_k &:= \text{vec}((\mathbf{a}'(\mathbf{t}_k) \odot \mathbf{r}_k) \mathbf{s}_{k-1}^\top) \\ \delta_k &:= \mathbf{r}_k \odot \mathbf{a}''(\mathbf{t}_k) + \mathbf{d}_k \odot \mathbf{a}'(\mathbf{t}_k)^2 \\ \mathbf{d}_{k-1} &:= \left(\sum_{j=1}^{D_k} \mathbf{W}_{k,ij}^2 \delta_{k,j} \right)_{i=1}^{D_k} \\ \mathbf{h}_k &:= \text{vec}(\delta_k (\mathbf{s}_{k-1}^2)^\top)\end{aligned}$$

starting from $\mathbf{r}_K = \nabla \ell(\mathbf{s}_K)$ and, e.g., $\mathbf{d}_K = \text{diag}(\nabla^2 \ell(\mathbf{s}_K))$. The algorithm returns $\mathbf{g}_1, \dots, \mathbf{g}_K$ as the gradients of f w.r.t. $\mathbf{w}_1, \dots, \mathbf{w}_K$, with $\mathbf{w}_i = \text{vec}(\mathbf{W}_i)$, and $\mathbf{h}_1, \dots, \mathbf{h}_K$ as the diagonal approximations of the Hessian w.r.t. $\mathbf{w}_1, \dots, \mathbf{w}_K$.

9.8 Randomized estimators

In this section, we describe randomized estimators of the diagonal of the Hessian or Gauss-Newton matrices.

9.8.1 Girard-Hutchinson estimator

We begin with a generic estimator, originally proposed for trace estimation by Girard (1989) and extended by Hutchinson (1989). Let $\mathbf{A} \in \mathbb{R}^{P \times P}$ be an arbitrary square matrix, whose matrix-vector product (matvec) is available. Suppose $\boldsymbol{\omega} \in \mathbb{R}^P$ is an isotropic random vector, i.e., such that $\mathbb{E}_{\boldsymbol{\omega} \sim p}[\boldsymbol{\omega} \boldsymbol{\omega}^\top] = \mathbf{I}$. For example, two common choices are the Rademacher distribution $p = \text{Uniform}(\{-1, 1\})$ and the standard normal distribution $p = \text{Normal}(0, \mathbf{I})$. Then, we have

$$\mathbb{E}_{\boldsymbol{\omega} \sim p}[\langle \boldsymbol{\omega}, \mathbf{A} \boldsymbol{\omega} \rangle] = \text{tr}(\mathbf{A}).$$

Applications include generalized cross-validation, computing the Kullback-Leibler divergence between two Gaussians, and computing the derivatives of the log-determinant.

The approach can be extended (Bekas *et al.*, 2007; Baston and Nakatsukasa, 2022; Hallman *et al.*, 2023) to obtain an estimator of the

diagonal of \mathbf{A} ,

$$\mathbb{E}_{\omega \sim p}[\omega \odot \mathbf{A}\omega] = \text{Diag}(\mathbf{A}),$$

where \odot denotes the Hadamard product (element-wise multiplication). This suggests that we can use the Monte-Carlo method to estimate the diagonal of \mathbf{A} ,

$$\text{Diag}(\mathbf{A}) \approx \frac{1}{S} \sum_{i=1}^S \omega_i \odot \mathbf{A}\omega_i,$$

with equality as $S \rightarrow \infty$, since the estimator is unbiased. Since, as reviewed in Section 9.1 and Section 9.2, we know how to multiply efficiently with the Hessian and the Gauss-Newton matrices, we can apply the technique with these matrices. The variance is determined by the number S of samples drawn and therefore by the number of matvecs performed. More elaborated approaches have been proposed to further reduce the variance (Meyer *et al.*, 2021; Epperly *et al.*, 2023).

9.8.2 Bartlett estimator for the factorization

Suppose the objective function is of the form $L(\mathbf{w}; \mathbf{x}, \mathbf{y}) := \ell(f(\mathbf{w}; \mathbf{x}); \mathbf{y})$ where ℓ is the negative log-likelihood $\ell(\boldsymbol{\theta}; \mathbf{y}) := -\log p_{\boldsymbol{\theta}}(\mathbf{y})$ of an exponential family distribution, and $\boldsymbol{\theta} := f(\mathbf{w}; \mathbf{x})$, for some network f . We saw from the equivalence between the Fisher and Gauss-Newton matrices in Proposition 9.6 (which follows from the Bartlett identity) that

$$\begin{aligned} \nabla_{\text{GN}}^2 L(\mathbf{w}; \mathbf{x}, \cdot) &= \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}} [\partial f(\mathbf{w}; \mathbf{x})^* \nabla \ell(\boldsymbol{\theta}; Y) \otimes \nabla \ell(\boldsymbol{\theta}; Y) \partial f(\mathbf{w}; \mathbf{x})] \\ &= \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}} [\nabla L(\mathbf{w}; \mathbf{x}, Y) \otimes \nabla L(\mathbf{w}; \mathbf{x}, Y)], \end{aligned}$$

where \cdot indicates that the result holds for any value of the second argument. This suggests a Monte-Carlo scheme

$$\nabla_{\text{GN}}^2 L(\mathbf{w}; \mathbf{x}, \cdot) \approx \frac{1}{S} \sum_{j=1}^S [\nabla L(\mathbf{w}; \mathbf{x}, \mathbf{y}_{i_j}) \otimes \nabla L(\mathbf{w}; \mathbf{x}, \mathbf{y}_{i_j})]$$

where $\mathbf{y}_{i_1}, \dots, \mathbf{y}_{i_S} \sim p_{\boldsymbol{\theta}}$ and $\boldsymbol{\theta} = f(\mathbf{w}, \mathbf{x})$. In words, we can approximate the Gauss-Newton matrix with S gradient computations. This factorization can also be used to approximate the GNVP in Eq. (9.1).

9.8.3 Bartlett estimator for the diagonal

Following a similar approach, we obtain

$$\mathbf{diag}(\nabla_{\text{GN}}^2 L(\mathbf{w}; \mathbf{x}, \cdot)) = \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}} [\nabla L(\mathbf{w}; \mathbf{x}, Y) \odot \nabla L(\mathbf{w}; \mathbf{x}, Y)],$$

where \odot indicates the element-wise (Hadamard) product. Using a Monte-Carlo scheme, sampling $\mathbf{y}_{i_1}, \dots, \mathbf{y}_{i_S}$ from $p_{\boldsymbol{\theta}}$, we therefore obtain

$$\mathbf{diag}(\nabla_{\text{GN}}^2 L(\mathbf{w}; \mathbf{x}, \cdot)) \approx \frac{1}{S} \sum_{j=1}^S \nabla L(\mathbf{w}; \mathbf{x}, \mathbf{y}_{i_j}) \odot \nabla L(\mathbf{w}; \mathbf{x}, \mathbf{y}_{i_j}),$$

with equality when all labels in the support of $p_{\boldsymbol{\theta}}$ have been sampled. That estimator, used for instance in (Wei *et al.*, 2020, Appendix C.1.), requires access to **individual gradients** evaluated at the sampled labels. Another possible estimator of the diagonal is given by

$$\begin{aligned} & \frac{1}{S} \mathbf{diag}(\nabla_{\text{GN}}^2 L(\mathbf{w}; \mathbf{x}, \cdot)) \\ &= \mathbb{E}_{Y_1, \dots, Y_S \sim p_{\boldsymbol{\theta}}} \left[\nabla \frac{1}{S} \sum_{i=1}^S L(\mathbf{w}; \mathbf{x}, Y_i) \odot \nabla \frac{1}{S} \sum_{i=1}^S L(\mathbf{w}; \mathbf{x}, Y_i) \right]. \end{aligned}$$

Letting $\gamma_i := \nabla L(\mathbf{w}; \mathbf{x}, Y_i)$, this follows from

$$\begin{aligned} \mathbb{E} \left[\sum_i \gamma_i \odot \sum_j \gamma_j \right] &= \mathbb{E} \left[\sum_i \gamma_i \odot \gamma_i + \sum_{i \neq j} \gamma_i \odot \gamma_j \right] \\ &= \mathbb{E} \left[\sum_i \gamma_i \odot \gamma_i \right] \end{aligned}$$

where we used that $\mathbb{E}[\gamma_i \odot \gamma_j] = \mathbb{E}[\gamma_i] \odot \mathbb{E}[\gamma_j] = \mathbf{0}$ since γ_i and γ_j are independent variables for $i \neq j$ and have zero mean, from Bartlett's first identity Eq. (12.2). We can then use the Monte-Carlo method to obtain

$$\frac{1}{S} \mathbf{diag}(\nabla_{\text{GN}}^2 L(\mathbf{w}; \mathbf{x}, \cdot)) \approx \left(\nabla \frac{1}{S} \sum_{j=1}^S L(\mathbf{w}; \mathbf{x}, \mathbf{y}_{i_j}) \right) \odot \left(\nabla \frac{1}{S} \sum_{j=1}^S L(\mathbf{w}; \mathbf{x}, \mathbf{y}_{i_j}) \right),$$

with equality when all labels in the support of $p_{\boldsymbol{\theta}}$ have been sampled. This estimator can be more convenient to implement, since it only needs access to the gradient of the **averaged** losses. However, it may suffer from higher variance. A special case of this estimator is used by Liu *et al.* (2023), where they draw only one \mathbf{y} for each \mathbf{x} .

9.9 Summary

- By using a Hessian chain rule, we can develop a “Hessian backpropagation”. While it is reasonably simple for computation chains, it becomes computationally prohibitive for computation graphs, due to the cross-product terms occurring with fan-in.
- A better approach is to use Hessian-vector products (HVPs). We saw that there are four possible methods to compute HVPs, but the forward-over-reverse method is a priori the most efficient. Similarly as for computing gradients, computing HVPs is only a constant times more expensive than evaluating the function itself.
- The Gauss-Newton matrix associated with the composition $\ell \circ f$ can be seen as an approximation of the Hessian. It is a positive semidefinite matrix if ℓ is convex, and can be used to build a principled quadratic approximation of a function. It is equivalent to the Fisher information matrix in the case of exponential families. Gauss-Newton-vector products can be computed efficiently, like HVPs.
- We also described other approximations, such as (block) diagonal approximations, and randomized estimators.

10

Inference in graphical models as differentiation

A graphical model specifies how random variables depend on each other and therefore determines how their joint probability distribution factorizes. In this chapter, we review key concepts in graphical models and how they relate to differentiation, drawing in the process analogies with computation chains and computation graphs.

10.1 Chain rule of probability

The chain rule of probability is a fundamental law in probability theory for computing the **joint probability** of events. In the case of only two events A_1 and A_2 , it reduces to the **product rule**

$$\mathbb{P}(A_1 \cap A_2) = \mathbb{P}(A_2|A_1)\mathbb{P}(A_1).$$

For two discrete random variables S_1 and S_2 , using the events $A_1 := \{S_1 = \mathbf{s}_1\}$ and $A_2 := \{S_2 = \mathbf{s}_2\}$, the product rule becomes

$$\mathbb{P}(S_1 = \mathbf{s}_1, S_2 = \mathbf{s}_2) = \mathbb{P}(S_2 = \mathbf{s}_2|S_1 = \mathbf{s}_1)\mathbb{P}(S_1 = \mathbf{s}_1).$$

More generally, using the product rule, we have for K events

$$\mathbb{P}(A_1 \cap \dots \cap A_K) = \mathbb{P}(A_K \mid A_1 \cap \dots \cap A_{K-1}) \mathbb{P}(A_1 \cap \dots \cap A_{K-1}).$$

Applying the product rule one more time, we have

$$\mathbb{P}(A_1 \cap \dots \cap A_{K-1}) = \mathbb{P}(A_{K-1} \mid A_1 \cap \dots \cap A_{K-2}) \mathbb{P}(A_1 \cap \dots \cap A_{K-2}).$$

Repeating the process recursively, we arrive at the **chain rule of probability**

$$\begin{aligned} \mathbb{P}(A_1 \cap \dots \cap A_K) &= \prod_{j=1}^K \mathbb{P}(A_j \mid A_1 \cap \dots \cap A_{j-1}) \\ &= \prod_{j=1}^K \mathbb{P}\left(A_j \mid \bigcap_{i=1}^{j-1} A_i\right). \end{aligned}$$

For K discrete random variables S_j , using the events $A_j := \{S_j = \mathbf{s}_j\}$, the chain rule of probability becomes

$$\mathbb{P}(S_1 = \mathbf{s}_1, \dots, S_K = \mathbf{s}_K) = \prod_{j=1}^K \mathbb{P}(S_j = \mathbf{s}_j \mid S_1 = \mathbf{s}_1, \dots, S_{j-1} = \mathbf{s}_{j-1}).$$

Importantly, this factorization holds **without** any independence assumption on the variables S_1, \dots, S_K . In other words, the space of probability distributions specified by the joint probability on the left-hand side, and the space of probability distributions specified by the product of conditional probabilities on the right-hand side, are the same. We can further simplify the factorization if we make additional conditional independence assumptions.

10.2 Conditional independence

We know that if two events A and B are independent, then

$$\mathbb{P}(A|B) = \mathbb{P}(A).$$

Similarly, if two random variables S_1 and S_2 are independent, then

$$\mathbb{P}(S_2 = \mathbf{s}_2 | S_1 = \mathbf{s}_1) = \mathbb{P}(S_2 = \mathbf{s}_2).$$

More generally, if we work with K variables S_1, \dots, S_K , some variables may depend on each other, while others may not. To simplify the

notation, given a set \mathcal{C} , we define the shorthands

$$\begin{aligned} S_{\mathcal{C}} &:= (S_i : i \in \mathcal{C}) \\ \mathbf{s}_{\mathcal{C}} &:= (\mathbf{s}_i : i \in \mathcal{C}). \end{aligned}$$

We say that a variable S_j is independent of $S_{\mathcal{D}}$ **conditioned** on $S_{\mathcal{C}}$, with $\mathcal{C} \cap \mathcal{D} = \emptyset$, if for any $\mathbf{s}_j, \mathbf{s}_{\mathcal{C}}, \mathbf{s}_{\mathcal{D}}$

$$\mathbb{P}(S_j = \mathbf{s}_j \mid S_{\mathcal{C}} = \mathbf{s}_{\mathcal{C}}, S_{\mathcal{D}} = \mathbf{s}_{\mathcal{D}}) = \mathbb{P}(S_j = \mathbf{s}_j \mid S_{\mathcal{C}} = \mathbf{s}_{\mathcal{C}}).$$

10.3 Inference problems

10.3.1 Joint probability distributions

We consider a collection of K variables $\mathbf{s} := (\mathbf{s}_1, \dots, \mathbf{s}_K)$, potentially **ordered** or **unordered**. Each \mathbf{s} belongs to the **Cartesian product** $\mathcal{S} := \mathcal{S}_1 \times \dots \times \mathcal{S}_K$. Throughout this chapter, we assume that the sets \mathcal{S}_k are discrete for concreteness, with $\mathcal{S}_k := \{\mathbf{v}_1, \dots, \mathbf{v}_{M_k}\}$. Note that because \mathcal{S}_k is discrete, we can always identify it with $\{1, \dots, M_k\}$. A graphical model specifies a **joint probability distribution**

$$\begin{aligned} \mathbb{P}(S = \mathbf{s}) &= \mathbb{P}(S_1 = \mathbf{s}_1, \dots, S_K = \mathbf{s}_K) \\ &= p(\mathbf{s}) \\ &= p(\mathbf{s}_1, \dots, \mathbf{s}_K), \end{aligned}$$

where p is the probability mass function of the joint probability distribution. Summing over the Cartesian product of all possible configurations, we obtain

$$\sum_{\mathbf{s} \in \mathcal{S}} p(\mathbf{s}) = \sum_{\mathbf{s}_1, \dots, \mathbf{s}_K \in \mathcal{S}} p(\mathbf{s}_1, \dots, \mathbf{s}_K) = 1.$$

As we shall see, the graph of a graphical model encodes the **dependencies** between the variables (S_1, \dots, S_K) and therefore how their joint distribution **factorizes**. Given access to a joint probability distribution, there are several **inference** problems one typically needs to solve.

10.3.2 Likelihood

A simple task is to compute the **likelihood** of some observations $\mathbf{s} = (\mathbf{s}_1, \dots, \mathbf{s}_K)$,

$$\mathbb{P}(S_1 = \mathbf{s}_1, \dots, S_k = \mathbf{s}_k) = p(\mathbf{s}_1, \dots, \mathbf{s}_k).$$

It is also common to compute the **log-likelihood**,

$$\log \mathbb{P}(S_1 = \mathbf{s}_1, \dots, S_k = \mathbf{s}_k) = \log p(\mathbf{s}_1, \dots, \mathbf{s}_k).$$

10.3.3 Maximum a-posteriori inference

Another common task is to compute the most likely configuration,

$$\arg \max_{\mathbf{s}_1 \in \mathcal{S}_1, \dots, \mathbf{s}_K \in \mathcal{S}_K} p(\mathbf{s}_1, \dots, \mathbf{s}_K).$$

This is the **mode** of the joint probability distribution. This is also known as maximum a-posteriori (MAP) inference in the literature (Wainwright and Jordan, 2008).

10.3.4 Marginal inference

The operation of **marginalization** consists in summing (or integrating) over all possible values of a given variable in a joint probability distribution. This allows us to compute the **marginal probability** of the remaining variables. For instance, we may want to marginalize all variables but $S_k = \mathbf{s}_k$. To do so, we define the Cartesian product

$$\mathcal{C}_k(\mathbf{s}_k) := \underbrace{\mathcal{S}_1 \times \dots \times \mathcal{S}_{k-1}}_{\mathcal{A}_{k-1}} \times \{\mathbf{s}_k\} \times \underbrace{\mathcal{S}_{k+1} \times \dots \times \mathcal{S}_K}_{\mathcal{B}_{k+1}}. \quad (10.1)$$

Summing over all variables but S_k , we obtain the marginal probability of $S_k = \mathbf{s}_k$ as

$$\begin{aligned} \mathbb{P}(S_k = \mathbf{s}_k) &= \sum_{\mathbf{s}_1, \dots, \mathbf{s}_K \in \mathcal{C}_k(\mathbf{s}_k)} p(\mathbf{s}_1, \dots, \mathbf{s}_K) \\ &= \sum_{\mathbf{s}_1, \dots, \mathbf{s}_{k-1} \in \mathcal{A}_{k-1}} \sum_{\mathbf{s}_{k+1}, \dots, \mathbf{s}_K \in \mathcal{B}_{k+1}} p(\mathbf{s}_1, \dots, \mathbf{s}_K) \end{aligned}$$

Defining similarly

$$\mathcal{C}_{k,l}(\mathbf{s}_k, \mathbf{s}_l) := \mathcal{S}_1 \times \dots \times \{\mathbf{s}_k\} \times \dots \times \{\mathbf{s}_l\} \times \dots \times \mathcal{S}_K,$$

we obtain

$$\mathbb{P}(S_k = \mathbf{s}_k, S_l = \mathbf{s}_l) = \sum_{\mathbf{s}_1, \dots, \mathbf{s}_K \in \mathcal{C}_{k,l}(\mathbf{s}_k, \mathbf{s}_l)} p(\mathbf{s}_1, \dots, \mathbf{s}_K).$$

In particular, we may want to compute the marginal probability of two consecutive variables, $\mathbb{P}(S_{k-1} = \mathbf{s}_{k-1}, S_k = \mathbf{s}_k)$.

10.3.5 Expectation, convex hull, marginal polytope

Another common operation is to compute the expectation of $\phi(S)$ under a distribution p . It is defined by

$$\boldsymbol{\mu} := \mathbb{E}_{S \sim p}[\phi(S)] = \sum_{s \in \mathcal{S}} p(s) \phi(s) \in \mathcal{M}$$

For the expectation under p_θ , we write

$$\mu(\theta) := \mathbb{E}_{S \sim p_\theta}[\phi(S)] = \sum_{s \in \mathcal{S}} p_\theta(s) \phi(s) \in \mathcal{M}.$$

In exponential family distributions (Section 3.4), the function ϕ is called a **statistic**. It decomposes as

$$\phi(s) := (\phi_{\mathcal{C}}(s_{\mathcal{C}}))_{\mathcal{C} \in \mathcal{C}},$$

where $\mathcal{C} \subseteq [K]$. Intuitively, $\phi(s)$ can be thought as an **encoding** or **embedding** of s (a potentially discrete object such as a sequence of integers) in a vector space. Under this decomposition, we can also compute

$$\mu_{\mathcal{C}} := \mathbb{E}_S[\phi_{\mathcal{C}}(S_{\mathcal{C}})] = \sum_{s \in \mathcal{S}} p(s) \phi_{\mathcal{C}}(s_{\mathcal{C}}).$$

Convex hull

The mean $\boldsymbol{\mu}$ belongs to the **convex hull** of $\phi(\mathcal{S}) := \{\phi(s) : s \in \mathcal{S}\}$,

$$\mathcal{M} := \text{conv}(\phi(\mathcal{S})) := \left\{ \sum_{s \in \mathcal{S}} p(s) \phi(s) : p \in \mathcal{P}(\mathcal{S}) \right\},$$

where $\mathcal{P}(\mathcal{S})$ is the set of all possible probability distributions over \mathcal{S} . In other words, \mathcal{M} is the set of all possible convex combinations of $\phi(s)$ for $s \in \mathcal{S}$. The vertices of \mathcal{M} are all the $s \in \mathcal{S}$.

Case of binary encodings: the marginal polytope

In the special case of a discrete set $\mathcal{S}_k = \{v_1, \dots, v_M\}$ and of a **binary encoding** (indicator function) $\phi(s)$, the set \mathcal{M} is called the **marginal polytope** (Wainwright and Jordan, 2008), because each point $\boldsymbol{\mu} \in$

\mathcal{M} contains marginal probabilities. To see why, consider the **unary** potential

$$[\phi(\mathbf{s})]_{k,i} = [\phi_k(\mathbf{s}_k)]_i = \mathbb{I}(\mathbf{s}_k = \mathbf{v}_i) \quad (10.2)$$

where $\mathbb{I}(p) := 1$ if p is true, 0 otherwise. We then obtain the marginal probability of $S_k = \mathbf{v}_i$,

$$\begin{aligned} [\boldsymbol{\mu}]_{k,i} &= \mathbb{E}_S[\phi(S)_{k,i}] \\ &= \mathbb{E}_{S_k}[\phi_k(S_k)_i] \\ &= \mathbb{E}_{S_k}[\mathbb{I}(S_k = \mathbf{v}_i)] \\ &= \sum_{\mathbf{s}_k \in \mathcal{S}_k} \mathbb{P}(S_k = \mathbf{s}_k) \mathbb{I}(\mathbf{s}_k = \mathbf{v}_i) \\ &= \mathbb{P}(S_k = \mathbf{v}_i). \end{aligned}$$

Likewise, consider the **pairwise** potential

$$[\phi(\mathbf{s})]_{k,l,i,j} = [\phi_{k,l}(\mathbf{s}_k, \mathbf{s}_l)]_{i,j} = \mathbb{I}(\mathbf{s}_k = \mathbf{v}_i, \mathbf{s}_l = \mathbf{v}_j). \quad (10.3)$$

We then obtain the marginal probability of $S_k = \mathbf{v}_i$ and $S_l = \mathbf{v}_j$,

$$\begin{aligned} [\boldsymbol{\mu}]_{k,l,i,j} &= \mathbb{E}_S[\phi(S)_{k,l,i,j}] \\ &= \mathbb{E}_{S_k, S_l}[\phi_{k,l}(S_k, S_l)_{i,j}] \\ &= \mathbb{E}_{S_k, S_l}[\mathbb{I}(S_k = \mathbf{v}_i, S_l = \mathbf{v}_j)] \\ &= \sum_{\mathbf{s}_k \in \mathcal{S}_k} \sum_{\mathbf{s}_l \in \mathcal{S}_l} \mathbb{P}(S_k = \mathbf{s}_k, S_l = \mathbf{s}_l) \mathbb{I}(\mathbf{s}_k = \mathbf{v}_i, \mathbf{s}_l = \mathbf{v}_j) \\ &= \mathbb{P}(S_k = \mathbf{v}_i, S_l = \mathbf{v}_j). \end{aligned}$$

We can do the same with higher-order potential functions.

10.3.6 Complexity of brute force

Apart from computing the likelihood, which is trivial, computing the marginal, mode and expectation by brute force takes $O(\prod_{k=1}^K |\mathcal{S}_k|)$ time. In particular, if $|\mathcal{S}_k| = M \ \forall k \in [K]$, brute force takes $O(M^K)$ time.

10.4 Markov chains

In this section, we briefly review Markov chains. Our notation is chosen to emphasize the analogies with computation chains.

10.4.1 The Markov property

When random variables are organized **sequentially** as S_1, \dots, S_K , a simple example of conditional independence is when each variable $S_k \in \mathcal{S}_k$ only depends on the previous variable $S_{k-1} \in \mathcal{S}_{k-1}$, that is,

$$\begin{aligned} \mathbb{P}(S_k = \mathbf{s}_k \mid S_{k-1} = \mathbf{s}_{k-1}, \dots, S_1 = \mathbf{s}_1) &= \mathbb{P}(S_k = \mathbf{s}_k \mid S_{k-1} = \mathbf{s}_{k-1}) \\ &:= p_k(\mathbf{s}_k \mid \mathbf{s}_{k-1}), \end{aligned}$$

A probability distribution satisfying the above is said to satisfy the **Markov property**, and is called a **Markov chain**. A computation chain is specified by the functions f_k , that take \mathbf{s}_{k-1} as input and output \mathbf{s}_k . In analogy, a Markov chain is specified by the **conditional** probability distributions p_k of S_k given S_{k-1} . We can then define the **generative process**

$$\begin{aligned} S_0 &:= \mathbf{s}_0 \\ S_1 &\sim p_1(\cdot \mid S_0) \\ S_2 &\sim p_2(\cdot \mid S_1) \\ &\vdots \\ S_K &\sim p_K(\cdot \mid S_{K-1}). \end{aligned}$$

Strictly speaking, we should write $S_k \mid S_{k-1} \sim p_k(\cdot \mid S_{k-1})$. We choose our notation both for conciseness and for analogy with computation chains. Furthermore, to simplify the notation, we assume without loss of generality that S_0 is deterministic (if this is not the case, we can always move S_0 to S_1 and add a dummy variable as S_0). That is, $\mathbb{P}(S_0 = \mathbf{s}_0) = p_0(\mathbf{s}_0) := 1$ and $\mathcal{S}_0 := \{\mathbf{s}_0\}$. This amounts to setting the **initial distribution** of S_1 as

$$\mathbb{P}(S_1 = \mathbf{s}_1) := \mathbb{P}(S_0 = \mathbf{s}_0)\mathbb{P}(S_1 = \mathbf{s}_1 \mid S_0 = \mathbf{s}_0) = \mathbb{P}(S_1 = \mathbf{s}_1 \mid S_0 = \mathbf{s}_0).$$

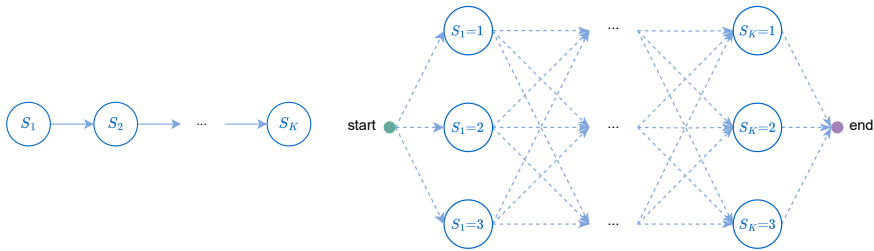


Figure 10.1: Left: Markov chain. Right: Computation graph of the forward-backward and the Viterbi algorithms: a lattice.

We can then compute the joint probability of the Markov chain by

$$\begin{aligned}
 \mathbb{P}(S_1 = \mathbf{s}_1, \dots, S_K = \mathbf{s}_K) &= p(\mathbf{s}_1, \dots, \mathbf{s}_K) \\
 &= \prod_{k=1}^K \mathbb{P}(S_k = \mathbf{s}_k \mid S_{k-1} = \mathbf{s}_{k-1}) \\
 &= \prod_{k=1}^K p_k(\mathbf{s}_k \mid \mathbf{s}_{k-1}),
 \end{aligned}$$

where we left the dependence on \mathbf{s}_0 implicit, since $p_0(\mathbf{s}_0) = 1$. A Markov chain with $S_k = \{1, 2, 3\}$ is illustrated in Fig. 10.1. A chain defines a **totally ordered set** $\{1, \dots, K\}$, since two nodes in the graph are necessarily linked to each other by a path.

Example 10.1 (Chain of categorical distributions). Suppose our goal is predict, from $\mathbf{x} \in \mathcal{X}$, a sequence of length K , where each S_k belongs to $\mathcal{S}_k = \{1, \dots, M\}$. In natural language processing, this task is called sequence tagging. We can define

$$S_k \sim \text{Categorical}(\boldsymbol{\pi}_{k-1,k,S_{k-1}})$$

where

$$\begin{aligned}
 \boldsymbol{\pi}_{k-1,k,i} &:= \text{softargmax}(\boldsymbol{\theta}_{k-1,k,i}) \in \triangle^M \\
 &= (\pi_{k-1,k,i,j})_{j=1}^M \\
 \boldsymbol{\theta}_{k-1,k,i} &:= (\theta_{k-1,k,i,j})_{j=1}^M \in \mathbb{R}^M \\
 \theta_{k-1,k,i,j} &:= f_{k-1,k}(\mathbf{x}, i, j, \mathbf{w}_k) \in \mathbb{R}.
 \end{aligned}$$

We therefore have

$$\begin{aligned}
 \mathbb{P}(S_k = j \mid S_{k-1} = i) &= p_k(j \mid i) \\
 &= \pi_{k-1,k,i,j} \\
 &= [\text{softargmax}(\boldsymbol{\theta}_{k-1,k,i})]_j \\
 &= \frac{\exp(\theta_{k-1,k,i,j})}{\sum_{j'} \exp(\theta_{k-1,k,i,j'})}
 \end{aligned}$$

and

$$\begin{aligned}
 \log \mathbb{P}(S_k = j \mid S_{k-1} = i) &= \log p_k(j \mid i) \\
 &= \theta_{k-1,k,i,j} - \text{logsumexp}(\boldsymbol{\theta}_{k-1,k,i}) \\
 &= \theta_{k-1,k,i,j} - \log \sum_{j'} \exp(\theta_{k-1,k,i,j'}).
 \end{aligned}$$

We emphasize that because $k-1$ and k are always consecutive, the representation $\theta_{k-1,k,i,j}$ is inefficient; we could use $\theta_{k,i,j}$ instead. Our notation is designed for consistency with Markov random fields.

10.4.2 Time-homogeneous Markov chains

A **time-homogeneous** discrete-time Markov chain corresponds to the case when the distribution of S_k given S_{k-1} is the same regardless of k :

$$p_1 = \cdots = p_K = p.$$

The **finite-space** case corresponds to when each $S_k \in \mathcal{S}$ can take a finite set of values $\mathcal{S} = \{\mathbf{v}_1, \dots, \mathbf{v}_M\}$ and

$$\mathbb{P}(S_k = \mathbf{v}_j \mid S_{k-1} = \mathbf{v}_i) = p(\mathbf{v}_j \mid \mathbf{v}_i) = \pi_{i,j},$$

where $\pi_{i,j} \in [0, 1]$ is the **transition probability** from \mathbf{v}_i to \mathbf{v}_j . Because the set $\mathcal{S} = \{\mathbf{v}_1, \dots, \mathbf{v}_M\}$ is discrete, we can always identify it with $\{1, \dots, M\}$. That is, we can instead write

$$\mathbb{P}(S_k = j \mid S_{k-1} = i) = p(j \mid i) = \pi_{i,j}.$$

10.4.3 Higher-order Markov chains

More generally, a n^{th} -order Markov chain may depend, not only on the last variable, but on the last n variables,

$$\begin{aligned} & \mathbb{P}(S_k = \mathbf{s}_k \mid S_{k-1} = \mathbf{s}_{k-1}, \dots, S_1 = \mathbf{s}_1) \\ &= \mathbb{P}(S_k = \mathbf{s}_k \mid S_{k-1} = \mathbf{s}_{k-1}, \dots, S_{k-n} = \mathbf{s}_{k-n}) \\ &= p_k(\mathbf{s}_k \mid \mathbf{s}_{k-1}, \dots, \mathbf{s}_{k-n}). \end{aligned}$$

Autoregressive models such as Transformers (Section 4.8) can be seen as specifying a higher-order Markov chain, with a context window of size n . The larger context makes exact inference using dynamic programming computationally intractable. This is why practitioners use **beam search** or **ancestral sampling** (Section 10.5.3) instead.

10.5 Bayesian networks

In this section, we briefly review Bayesian networks. Our notation is chosen to emphasize the analogies with computation graphs.

10.5.1 Expressing variable dependencies using DAGs

Markov chains and more generally higher-order Markov chains are a special case of Bayesian network. Similarly to computation graphs reviewed in Section 8.3, variable dependencies can be expressed using a directed acyclic graph (DAG) $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where the vertices $\mathcal{V} = \{1, \dots, K\}$ represent variables and edges \mathcal{E} represent variable dependencies. The set $\{i_1, \dots, i_{n_k}\} = \text{pa}(k) \subseteq \mathcal{V}$, where $n_k := |\text{pa}(k)|$, indicates the variables $S_{i_1}, \dots, S_{i_{n_k}}$ that S_k depends on. This defines a **partially ordered set** (poset). For notational simplicity, we again assume without loss of generality that S_0 is deterministic. A computation graph is specified by functions f_1, \dots, f_K in topological order. In analogy, a **Bayesian network** is specified by **conditional** probability distributions p_k of S_k

given $S_{\text{pa}(k)}$. We can then define the **generative process**

$$\begin{aligned} S_0 &:= \mathbf{s}_0 \\ S_1 &\sim p_1(\cdot \mid S_0) \\ S_2 &\sim p_2(\cdot \mid S_{\text{pa}(2)}) \\ &\vdots \\ S_K &\sim p_K(\cdot \mid S_{\text{pa}(K)}). \end{aligned}$$

Using the chain rule of probability and variable independencies expressed by the DAG, the **joint probability distribution** is then (assuming a topological order for S_0, S_1, \dots, S_K)

$$\begin{aligned} \mathbb{P}(S = \mathbf{s}) &:= \mathbb{P}(S_1 = \mathbf{s}_1, \dots, S_K = \mathbf{s}_K) \\ &= \prod_{k=1}^K \mathbb{P}(S_k = s_k \mid S_{\text{pa}(k)} = \mathbf{s}_{\text{pa}(k)}) \\ &:= \prod_{k=1}^K p_k(s_k \mid \mathbf{s}_{\text{pa}(k)}) \end{aligned}$$

This representation is well suited to express **causal** relationships between random variables.

10.5.2 Parameterizing Bayesian networks

In a Bayesian framework, observed data, latent variables, parameters and noise variables are all treated as random variables. If the conditional distribution p_k associated to node k depends on some parameters, they can be provided to p_k as conditioning, using parent nodes.

A Bayesian network is specified by the conditional distributions p_k . Therefore, unlike computation graphs, there is no notion of function f_k in a Bayesian network. However, the root nodes of the Bayesian network can be the output of a neural network. For instance, autoregressive models, such as RNNs or Transformers, specify the conditional probability distribution of a token given past tokens, and the chain rule of probability is used to obtain a probability distribution over entire sequences.

10.5.3 Ancestral sampling

A major advantage of Bayesian networks is that, provided that each conditional distribution p_k is normalized, the joint distribution of $S = (S_1, \dots, S_K)$ is automatically normalized. This means that we can very easily draw i.i.d. samples from the joint distribution, by following the generative process: we follow the topological order $k = 1, \dots, K$ and on iteration k we draw a value $s_k \sim p_k(\cdot | \mathbf{s}_{\text{pa}(k)})$ conditioned on the previous values $\mathbf{s}_{\text{pa}(k)}$. This is known as **ancestral sampling**.

10.6 Markov random fields

10.6.1 Expressing factors using undirected graphs

A Markov random field (MRF), a.k.a. undirected graphical model, specifies a distribution that factorizes as

$$\mathbb{P}(S = \mathbf{s}) = p(\mathbf{s}) := \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{s}_C),$$

where C is the set of maximal **cliques** of \mathcal{G} , that is, subsets of \mathcal{V} that are fully connected, Z is a normalization constant defined by

$$Z := \sum_{\mathbf{s} \in \mathcal{S}} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{s}_C),$$

and $\psi_C: \mathcal{S}_C \rightarrow \mathbb{R}_+$ is a **potential function** (a.k.a. compatibility function), with $\mathcal{S}_C := (\mathcal{S}_j)_{j \in C}$. According to the Hammersley-Clifford theorem, an MRF can be equivalently defined in terms of Markov properties; we refer the interested reader to Wainwright and Jordan (2008). For the sake of this chapter, the definition above is sufficient for our purposes.

Example 10.2 (Markov chains as Markov random fields). For a chain, letting $S = (S_1, \dots, S_K)$ and $\mathbf{s} = (\mathbf{s}_1, \dots, \mathbf{s}_K)$, recall that

$$\mathbb{P}(S = \mathbf{s}) = \prod_{k=1}^K p_k(\mathbf{s}_k | \mathbf{s}_{k-1}).$$

This is equivalent to an MRF with $Z = 1$ (since a chain is auto-

matically normalized),

$$C := \{\{0, 1\}, \{1, 2\}, \dots, \{K-1, K\}\}$$

and with potential function

$$\psi_{\{k-1, k\}}(\mathbf{s}_{k-1}, \mathbf{s}_k) := p_k(\mathbf{s}_k | \mathbf{s}_{k-1}).$$

More generally, a Bayesian network can be similarly written as an MRF by creating appropriate potential functions corresponding to the parents of each node.

10.6.2 MRFs as exponential family distributions

Let us define the potential functions

$$\psi_{\mathcal{C}}(\mathbf{s}_{\mathcal{C}}; \boldsymbol{\theta}_{\mathcal{C}}) := \exp(\langle \boldsymbol{\theta}_{\mathcal{C}}, \phi_{\mathcal{C}}(\mathbf{s}_{\mathcal{C}}) \rangle)$$

for some sufficient statistic function $\phi_{\mathcal{C}}: \mathcal{S}_{\mathcal{C}} \rightarrow \Theta_{\mathcal{C}}$ and parameters $\boldsymbol{\theta}_{\mathcal{C}} \in \Theta_{\mathcal{C}}$. Then,

$$\begin{aligned} p_{\boldsymbol{\theta}}(\mathbf{s}) &:= \frac{1}{Z(\boldsymbol{\theta})} \prod_{\mathcal{C} \in \mathcal{C}} \psi_{\mathcal{C}}(\mathbf{s}_{\mathcal{C}}; \boldsymbol{\theta}_{\mathcal{C}}) \\ &= \frac{1}{Z(\boldsymbol{\theta})} \prod_{\mathcal{C} \in \mathcal{C}} \exp(\langle \boldsymbol{\theta}_{\mathcal{C}}, \phi_{\mathcal{C}}(\mathbf{s}_{\mathcal{C}}) \rangle) \\ &= \frac{1}{Z(\boldsymbol{\theta})} \exp \left(\sum_{\mathcal{C} \in \mathcal{C}} \langle \boldsymbol{\theta}_{\mathcal{C}}, \phi_{\mathcal{C}}(\mathbf{s}_{\mathcal{C}}) \rangle \right) \\ &= \frac{1}{Z(\boldsymbol{\theta})} \exp(\langle \boldsymbol{\theta}, \phi(\mathbf{s}) \rangle) \\ &= \exp(\langle \boldsymbol{\theta}, \phi(\mathbf{s}) \rangle - A(\boldsymbol{\theta})) \end{aligned}$$

where

$$\begin{aligned}
 \phi(\mathbf{s}) &:= (\phi_C(\mathbf{s}_C))_{C \in \mathcal{C}} \\
 \boldsymbol{\theta} &:= (\boldsymbol{\theta}_C)_{C \in \mathcal{C}} \\
 Z(\boldsymbol{\theta}) &:= \sum_{\mathbf{s} \in \mathcal{S}} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{s}_C; \boldsymbol{\theta}_C) \\
 &= \sum_{\mathbf{s} \in \mathcal{S}} \exp(\langle \boldsymbol{\theta}, \phi(\mathbf{s}) \rangle) \\
 A(\boldsymbol{\theta}) &:= \log Z(\boldsymbol{\theta})
 \end{aligned}$$

Therefore, for this choice of potential functions, we can view an MRF as an exponential family distribution (Section 3.4) with **natural parameters** $\boldsymbol{\theta}$, **sufficient statistic** ϕ and **log-partition function** $A(\boldsymbol{\theta})$.

Example 10.3 (Ising model). The Ising model is a classical example of MRF. Let $Y = (Y_1, \dots, Y_M) \in \{0, 1\}^M$ be an unordered collection of binary variables $Y_i \in \{0, 1\}$. This forms a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = [M]$ and $\mathcal{E} \subseteq V^2$, such that $(i, j) \in \mathcal{E}$ means that Y_i interacts with Y_j . In statistical physics, Y_i may indicate the presence or absence of particles, or the orientation of magnets. In image processing, Y_i may represent a black and white pixel. In multi-label classification, Y_i may indicate the presence or absence of a label. The probability of $\mathbf{y} = (y_1, \dots, y_M) \in \{0, 1\}^M$ is then

$$\begin{aligned}
 \mathbb{P}(Y = \mathbf{y}) &= p_{\boldsymbol{\theta}}(\mathbf{y}) \\
 &= \exp \left(\sum_{i \in \mathcal{V}} \theta_i y_i + \sum_{(i,j) \in \mathcal{E}} \theta_{i,j} y_i y_j - A(\boldsymbol{\theta}) \right) \\
 &= \exp \left(\sum_{C \in \mathcal{C}} \langle \boldsymbol{\theta}_C, \phi_C(\mathbf{y}) \rangle - A(\boldsymbol{\theta}) \right),
 \end{aligned}$$

where $C := \mathcal{V} \cup \mathcal{E}$ and $\boldsymbol{\theta} \in \mathbb{R}^{|\mathcal{V}|+|\mathcal{E}|}$ is the concatenation of $(\theta_i)_{i \in \mathcal{V}}$ and $(\theta_{i,j})_{(i,j) \in \mathcal{E}}$. These models are also known as **Boltzmann machines** in a neural network context. MAP inference in general Ising models is known to be NP-hard, but when the interaction weights $\theta_{i,j}$ are non-negative, MAP inference can be reduced to

graph cut algorithms (Greig *et al.*, 1989). There are two ways the above equation can be extended. First, we can use higher-order interactions, such as $y_i y_j y_k$ for $(i, j, k) \in \mathcal{V}^3$. Second, we may want to use categorical variables, which leads to the **Potts model**.

10.6.3 Conditional random fields

Conditional random fields (Lafferty *et al.*, 2001; Sutton, McCallum, *et al.*, 2012) are a special case of Markov random field, in which a conditioning variable is **explicitly** incorporated. For example, when the goal is to predict a variable \mathbf{y} conditioned on a variable \mathbf{x} , CRFs are defined as

$$\mathbb{P}(Y = \mathbf{y} \mid X = \mathbf{x}) = p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{\mathcal{C} \in \mathcal{C}} \psi_{\mathcal{C}}(\mathbf{y}_{\mathcal{C}}, \mathbf{x}).$$

Note that the potential functions $\Psi_{\mathcal{C}}$ are allowed to depend on the whole \mathbf{x} , as \mathbf{x} is just a conditioning variable.

10.6.4 Sampling

Contrary to Bayesian networks, MRFs require an explicit normalization constant Z . As a result, sampling from a distribution represented by a general MRF is usually more involved than for Bayesian networks. A commonly-used technique is **Gibbs sampling**.

10.7 Inference on chains

In this section, we review how to perform marginal inference and maximum a-posteriori inference on joint distributions of the form

$$p(\mathbf{s}_1, \dots, \mathbf{s}_K) = \frac{1}{Z} \prod_{k=1}^K \psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k),$$

where

$$Z := \sum_{\mathbf{s} \in \mathcal{S}} \prod_{k=1}^K \psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k)$$

and where we used ψ_k as a shorthand for $\psi_{k-1,k}$, since $k-1$ and k are consecutive. As explained in Example 10.2, this also includes Markov

chains by setting

$$\psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k) := p_k(\mathbf{s}_k \mid \mathbf{s}_{k-1}),$$

in which case $Z = 1$.

10.7.1 The forward-backward algorithm

The key idea of the forward-backward algorithm is to use the **distributivity** of multiplication over addition to write

$$Z = \sum_{\mathbf{s}_1 \in \mathcal{S}_1} \psi_1(\mathbf{s}_0, \mathbf{s}_1) \sum_{\mathbf{s}_2 \in \mathcal{S}_2} \psi_2(\mathbf{s}_1, \mathbf{s}_2) \cdots \sum_{\mathbf{s}_K \in \mathcal{S}_K} \psi_K(\mathbf{s}_{K-1}, \mathbf{s}_K).$$

We can compute these sums recursively, either **forward** or **backward**. Recalling the definitions of \mathcal{A}_{k-1} and \mathcal{B}_{k+1} in Eq. (10.1), we define the summations **up to** and **down to** k ,

$$\begin{aligned} \alpha_k(\mathbf{s}_k) &:= \sum_{\mathbf{s}_1, \dots, \mathbf{s}_{k-1} \in \mathcal{A}_{k-1}} \prod_{j=1}^k \psi_j(\mathbf{s}_{j-1}, \mathbf{s}_j) \\ &= \sum_{\mathbf{s}_{k-1} \in \mathcal{S}_{k-1}} \psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k) \cdots \sum_{\mathbf{s}_1 \in \mathcal{S}_1} \psi_2(\mathbf{s}_1, \mathbf{s}_2) \psi_1(\mathbf{s}_0, \mathbf{s}_1) \\ \beta_k(\mathbf{s}_k) &:= \sum_{\mathbf{s}_{k+1}, \dots, \mathbf{s}_K \in \mathcal{B}_{k+1}} \prod_{j=k+1}^K \psi_j(\mathbf{s}_{j-1}, \mathbf{s}_j) \\ &= \sum_{\mathbf{s}_{k+1} \in \mathcal{S}_{k+1}} \psi_{k+1}(\mathbf{s}_k, \mathbf{s}_{k+1}) \cdots \sum_{\mathbf{s}_K \in \mathcal{S}_K} \psi_K(\mathbf{s}_{K-1}, \mathbf{s}_K). \end{aligned}$$

We can compute the two quantities by recursing forward and backward

$$\begin{aligned} \alpha_k(\mathbf{s}_k) &= \sum_{\mathbf{s}_{k-1} \in \mathcal{S}_{k-1}} \psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k) \alpha_{k-1}(\mathbf{s}_{k-1}) \\ \beta_k(\mathbf{s}_k) &= \sum_{\mathbf{s}_{k+1} \in \mathcal{S}_{k+1}} \psi_{k+1}(\mathbf{s}_k, \mathbf{s}_{k+1}) \beta_{k+1}(\mathbf{s}_{k+1}) \end{aligned}$$

where we defined the initializations

$$\begin{aligned} \alpha_1(\mathbf{s}_1) &:= \psi_1(\mathbf{s}_0, \mathbf{s}_1) \quad \forall \mathbf{s}_1 \in \mathcal{S}_1 \\ \beta_K(\mathbf{s}_K) &:= 1 \quad \forall \mathbf{s}_K \in \mathcal{S}_K. \end{aligned}$$

The **normalization** term can then be computed by

$$Z = \sum_{\mathbf{s}_K \in \mathcal{S}_K} \alpha_K(\mathbf{s}_K) \beta_K(\mathbf{s}_K) = \sum_{\mathbf{s}_1 \in \mathcal{S}_1} \alpha_1(\mathbf{s}_1) \beta_1(\mathbf{s}_1)$$

and the **marginal probabilities** by

$$\begin{aligned} \mathbb{P}(S_k = \mathbf{s}_k) &= \frac{1}{Z} \alpha_k(\mathbf{s}_k) \beta_k(\mathbf{s}_k) \\ \mathbb{P}(S_{k-1} = \mathbf{s}_{k-1}, S_k = \mathbf{s}_k) &= \frac{1}{Z} \alpha_{k-1}(\mathbf{s}_{k-1}) \psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k) \beta_k(\mathbf{s}_k). \end{aligned}$$

We can also compute the conditional probabilities by

$$\begin{aligned} \mathbb{P}(S_k = \mathbf{s}_k \mid S_{k-1} = \mathbf{s}_{k-1}) &= \frac{\mathbb{P}(S_{k-1} = \mathbf{s}_{k-1}, S_k = \mathbf{s}_k)}{\mathbb{P}(S_{k-1} = \mathbf{s}_{k-1})} \\ &= \frac{\alpha_{k-1}(\mathbf{s}_{k-1}) \psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k) \beta_k(\mathbf{s}_k)}{\alpha_{k-1}(\mathbf{s}_{k-1}) \beta_{k-1}(\mathbf{s}_{k-1})} \\ &= \frac{\psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k) \beta_k(\mathbf{s}_k)}{\beta_{k-1}(\mathbf{s}_{k-1})}. \end{aligned}$$

In practice, the two recursions are often implemented in the **log-domain** for numerical stability,

$$\begin{aligned} \log \alpha_k(\mathbf{s}_k) &= \log \sum_{\mathbf{s}_{k-1} \in \mathcal{S}_{k-1}} \exp(\log \psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k) + \log \alpha_{k-1}(\mathbf{s}_{k-1})) \\ \log \beta_k(\mathbf{s}_k) &= \log \sum_{\mathbf{s}_{k+1} \in \mathcal{S}_{k+1}} \exp(\log \psi_{k+1}(\mathbf{s}_k, \mathbf{s}_{k+1}) + \log \beta_{k+1}(\mathbf{s}_{k+1})). \end{aligned}$$

We recognize the **log-sum-exp** operator, which can be implemented in a numerically stable way (Section 4.4.2). The overall **dynamic programming** procedure, a.k.a. **forward-backward** algorithm (Baum and Petrie, 1966; Rabiner, 1989), is summarized in Algorithm 10.1. We notice that the forward and backward passes are actually independent of each other, and can therefore be performed in parallel.

10.7.2 The Viterbi algorithm

Similarly, using the distributivity of multiplication over maximization,

$$\begin{aligned} &\max_{\mathbf{s}_1 \in \mathcal{S}_1, \dots, \mathbf{s}_K \in \mathcal{S}_K} \prod_{k=1}^K \psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k) \\ &= \max_{\mathbf{s}_K \in \mathcal{S}_K} \max_{\mathbf{s}_{K-1} \in \mathcal{S}_{K-1}} \psi_K(\mathbf{s}_{K-1}, \mathbf{s}_K) \dots \max_{\mathbf{s}_1 \in \mathcal{S}_1} \psi_2(\mathbf{s}_1, \mathbf{s}_2) \psi_1(\mathbf{s}_0, \mathbf{s}_1). \end{aligned}$$

Algorithm 10.1 Marginal inference on a chain**Potential functions:** ψ_1, \dots, ψ_K **Input:** s_0

- 1: Initialize $\alpha_1(s_1) := \psi_1(s_0, s_1) \forall s_1 \in \mathcal{S}_1$
 - 2: **for** $k := 2, \dots, K$ **do** ▷ Forward pass
 - 3: **for** $s_k \in \mathcal{S}_k$ **do**
 - 4: $\alpha_k(s_k) := \sum_{s_{k-1} \in \mathcal{S}_{k-1}} \psi_k(s_{k-1}, s_k) \alpha_{k-1}(s_{k-1})$
 - 5: Initialize $\beta_K(s_K) := 1 \forall s_K \in \mathcal{S}_K$
 - 6: **for** $k := K - 1, \dots, 1$ **do** ▷ Backward pass
 - 7: **for** $s_k \in \mathcal{S}_k$ **do**
 - 8: $\beta_k(s_k) := \sum_{s_{k+1} \in \mathcal{S}_{k+1}} \psi_{k+1}(s_k, s_{k+1}) \beta_{k+1}(s_{k+1})$
 - 9: Compute $Z = \sum_{s_K \in \mathcal{S}_K} \alpha_K(s_K) \beta_K(s_K) = \sum_{s_K \in \mathcal{S}_K} \alpha_K(s_K)$
- Outputs:** $\forall k \in [K]$:
- $$\mathbb{P}(S_k = s_k) = \frac{1}{Z} \alpha_k(s_k) \beta_k(s_k)$$
- $$\mathbb{P}(S_{k-1} = s_{k-1}, S_k = s_k) = \frac{1}{Z} \alpha_{k-1}(s_{k-1}) \psi_k(s_{k-1}, s_k) \beta_k(s_k)$$

Let us define for $k \in [K]$

$$\delta_k(s_k) := \max_{s_{k-1} \in \mathcal{S}_{k-1}} \psi_k(s_{k-1}, s_k) \dots \max_{s_1 \in \mathcal{S}_1} \psi_2(s_1, s_2) \psi_1(s_0, s_1).$$

We can compute these quantities recursively, since for $k \in [K]$

$$\delta_k(s_k) = \max_{s_{k-1} \in \mathcal{S}_{k-1}} \psi_k(s_{k-1}, s_k) \delta_{k-1}(s_{k-1}),$$

with $\delta_1(s_k) := \psi(s_0, s_k)$. We finally have

$$\max_{s_1 \in \mathcal{S}_1, \dots, s_K \in \mathcal{S}_K} p(s_1, \dots, s_K) = \frac{1}{Z} \max_{s_K \in \mathcal{S}_K} \delta_K(s_K).$$

In practice, for numerical stability, we often implement the forward recursion in the **log-domain**. Using the fact that the logarithm is monotonic, we indeed have for all $k \in [K]$

$$\log \delta_k(s_k) = \max_{s_{k-1} \in \mathcal{S}_{k-1}} \log \psi_k(s_{k-1}, s_k) + \log \delta_{k-1}(s_{k-1}).$$

To enable efficient **backtracking**, during the forward pass, we compute

$$q_k(\mathbf{s}_k) := \arg \max_{\mathbf{s}_{k-1} \in \mathcal{S}_{k-1}} \psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k) \delta_{k-1}(\mathbf{s}_{k-1})$$

which can be thought as **backpointers** from \mathbf{s}_k^* to \mathbf{s}_{k-1}^* .

The resulting dynamic programming procedure, a.k.a. **Viterbi algorithm** (Viterbi, 1967; Forney, 1973), is summarized in Algorithm 10.2.

Algorithm 10.2 MAP inference on a chain

Potential functions: ψ_1, \dots, ψ_K

Input: \mathbf{s}_0

- 1: Initialize $\delta_1(\mathbf{s}_1) := \psi_1(\mathbf{s}_0, \mathbf{s}_1) \forall \mathbf{s}_1 \in \mathcal{S}_1$
- 2: **for** $k := 2, \dots, K$ **do** ▷ Forward pass
- 3: **for** $\mathbf{s}_k \in \mathcal{S}_k$ **do**
- 4: $\delta_k(\mathbf{s}_k) := \max_{\mathbf{s}_{k-1} \in \mathcal{S}_{k-1}} \psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k) \delta_{k-1}(\mathbf{s}_{k-1})$
- 5: $q_k(\mathbf{s}_k) := \arg \max_{\mathbf{s}_{k-1} \in \mathcal{S}_{k-1}} \psi_k(\mathbf{s}_{k-1}, \mathbf{s}_k) \delta_{k-1}(\mathbf{s}_{k-1})$
- 6: $\delta^* := \max_{\mathbf{s}_K \in \mathcal{S}_K} \delta_K(\mathbf{s}_K)$
- 7: $\mathbf{s}_K^* := \arg \max_{\mathbf{s}_K \in \mathcal{S}_K} \delta_K(\mathbf{s}_K)$
- 8: **for** $k := K - 1, \dots, 1$ **do** ▷ Backtracking
- 9: $\mathbf{s}_k^* := q_{k+1}(\mathbf{s}_{k+1}^*)$

Outputs: $\max_{\mathbf{s}_1 \in \mathcal{S}_1, \dots, \mathbf{s}_K \in \mathcal{S}_K} p(\mathbf{s}_1, \dots, \mathbf{s}_K) \propto \delta^*$
 $\arg \max_{\mathbf{s}_1 \in \mathcal{S}_1, \dots, \mathbf{s}_K \in \mathcal{S}_K} p(\mathbf{s}_1, \dots, \mathbf{s}_K) = (\mathbf{s}_1^*, \dots, \mathbf{s}_K^*)$

10.8 Inference on trees

More generally, efficient inference based on dynamic programming can be performed when dependencies between variables are expressed using a **tree** or **polytree**. The resulting marginal inference and MAP inference algorithms are often referred to as the **sum-product** and **max-sum** algorithms. The sum-product algorithm is also known as **belief propagation** or **message passing**, since it can be interpreted as propagating “local messages” through the graph. See for instance (Wainwright and Jordan, 2008, Section 2.5.1) for more details.

10.9 Inference as differentiation

In this section, we review the profound connections between differentiating the log-partition function of an exponential family distribution on one hand, and performing marginal inference (as well as maximum a-posteriori inference in the zero-temperature limit) on the other hand.

10.9.1 Inference as gradient of the log-partition

We first discuss a well-known fact in the graphical model literature: when using a binary encoding as the sufficient statistic ϕ in an exponential family distribution, the gradient $\nabla A(\boldsymbol{\theta})$ of the log-partition $A(\boldsymbol{\theta})$ gathers all the marginals (Wainwright and Jordan, 2008).

To see why, recall from Section 3.4 the definition of an exponential family distribution

$$p_{\boldsymbol{\theta}}(\mathbf{s}) = h(\mathbf{s}) \exp [\langle \boldsymbol{\theta}, \phi(\mathbf{s}) \rangle - A(\boldsymbol{\theta})]$$

and of its log-partition

$$A(\boldsymbol{\theta}) := \log \sum_{\mathbf{s} \in \mathcal{S}} h(\mathbf{s}) \exp [\langle \boldsymbol{\theta}, \phi(\mathbf{s}) \rangle].$$

From Proposition 3.2,

$$\mu(\boldsymbol{\theta}) := \nabla A(\boldsymbol{\theta}) = \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}}[\phi(Y)] \in \mathcal{M}.$$

Therefore, with the **binary encodings** in Eq. (10.2) and Eq. (10.3),

$$\begin{aligned} \mathbb{P}(S_k = \mathbf{v}_i) &= [\nabla A(\boldsymbol{\theta})]_{k,i} \\ \mathbb{P}(S_k = \mathbf{v}_i, S_l = \mathbf{v}_j) &= [\nabla A(\boldsymbol{\theta})]_{k,l,i,j}. \end{aligned}$$

Put differently, if we have an efficient algorithm for computing $A(\boldsymbol{\theta})$, we can perform **reverse-mode autodiff** on $A(\boldsymbol{\theta})$ to obtain $\nabla A(\boldsymbol{\theta})$, and therefore obtain the marginal probabilities. Following Section 8.3.3, the complexity of computing all marginal probabilities is therefore roughly the same as that of computing $A(\boldsymbol{\theta})$.

In the special case of chains, we obtain

$$\begin{aligned} \mathbb{P}(S_k = \mathbf{v}_i) &= [\nabla A(\boldsymbol{\theta})]_{k,i} = \frac{1}{Z} \alpha_k(\mathbf{v}_i) \beta_k(\mathbf{v}_i) \\ \mathbb{P}(S_{k-1} = \mathbf{v}_i, S_k = \mathbf{v}_j) &= [\nabla A(\boldsymbol{\theta})]_{k-1,k,i,j} = \frac{1}{Z} \alpha_{k-1}(\mathbf{v}_i) \psi_k(\mathbf{v}_i, \mathbf{v}_j) \beta_k(\mathbf{v}_j), \end{aligned}$$

where we left the dependence of Z , α and β on θ implicit.

If we define $A_\varepsilon(\theta) := \varepsilon A(\theta/\varepsilon)$, in the zero-temperature limit $\varepsilon \rightarrow 0$, we obtain that $\mu(\theta)$ is a binary encoding of the mode, i.e., of the maximum a-posteriori inference solution.

We now show i) how to unify the forward pass of the forward-backward and Viterbi algorithms using semirings and softmax operators
ii) how to compute the gradient of the log-partition using backpropagation.

10.9.2 Semirings and softmax operators

The forward passes in the forward-backward and Viterbi algorithms are clearly similar. In fact, they can be formally linked to each other using **semirings**.

Definition 10.1 (Semiring). A semiring is a set \mathbb{K} equipped with two binary operations (\oplus, \otimes) such that

- \otimes is commutative and associative,
- \oplus is associative and distributive over \otimes ,
- \otimes and \oplus have identity element $\bar{0}$ and $\bar{1}$, respectively.

We use the notations \oplus , \otimes , $\bar{0}$ and $\bar{1}$ to clearly distinguish them from the classical addition, multiplication, 0 and 1.

We recall the following laws for binary operations:

- **Commutativity** of \oplus : $a \oplus b = b \oplus a$,
- **Associativity** of \oplus : $a \oplus (b \oplus c) = (a \oplus b) \oplus c$,
- **Distributivity** of \otimes over \oplus : $a \otimes (b \oplus c) = (a \otimes b) \oplus (a \otimes c)$.

A set equipped with a binary operation supporting associativity and an identity element is called a **monoid**. A monoid such that every element has an inverse element is called a **group**. The difference between a ring and a semiring is that the latter only requires (\mathbb{K}, \oplus) and (\mathbb{K}, \otimes) to be monoids, not groups.

Equipped with these definitions, we can interpret the forward passes in the Viterbi and forward-backward algorithms as follows:

- the forward-backward algorithm in the exponential domain uses the semiring \mathbb{R}_+ equipped with $(+, \times)$ and identity elements $(0, 1)$;
- the Viterbi algorithm in the log domain uses the semiring \mathbb{R} equipped with $(\max, +)$ and identity elements $(-\infty, 0)$;
- the forward-backward algorithm in the log domain uses the semiring \mathbb{R} equipped with $(\max_\varepsilon, +)$ and identity elements $(-\infty, 0)$, where we defined the soft max operator (log-add-exp)

$$\max_\varepsilon(a, b) := \varepsilon \log((\exp(a) + \exp(b))/\varepsilon),$$

with $\varepsilon := 1$ by default.

It can be checked that indeed \max_ε is commutative, associative, and addition is distributive over \max_ε . Its identity element is $-\infty$. By associativity,

$$\begin{aligned} \max_\varepsilon(a_1, \max_\varepsilon(a_2, a_3)) &= \text{logsumexp}_\varepsilon(a_1, a_2, a_3) \\ &= \varepsilon \log \sum_i \exp(a_i/\varepsilon). \end{aligned}$$

In contrast, note that the sparsemax in Section 13.5 is not associative.

Thanks to associativity, we can introduce the shorthand notations

$$\max_{\varepsilon} f(\mathbf{v}) := \varepsilon \log \sum_{\mathbf{v} \in \mathcal{V}} \exp(f(\mathbf{v})/\varepsilon) \in \mathbb{R}.$$

and

$$\operatorname{argmax}_{\varepsilon} f(\mathbf{v}) := \left(\exp(f(\mathbf{v}')/\varepsilon) / \sum_{\mathbf{v} \in \mathcal{V}} \exp(f(\mathbf{v})/\varepsilon) \right)_{\mathbf{v}' \in \mathcal{V}} \in \mathcal{P}(\mathcal{V}).$$

Many algorithms can be generalized thanks to the use of semirings; see among others Aji and McEliece (2000) and Mohri *et al.* (2008). The distributive and associative properties play a key role in breaking down large problems into smaller ones (Verdu and Poor, 1987).

10.9.3 Inference as backpropagation

In this section, we show that, algorithmically, backtracking is recovered as a special case of backpropagation. See also (Eisner, 2016; Mensch and Blondel, 2018).

For notation simplicity, we assume $\mathcal{S}_0 = \{1\}$ and $\mathcal{S}_k = \{1, \dots, M\}$ for all $k \in [K]$. We focus on the case

$$\log \psi_k(i, j) = \langle \boldsymbol{\theta}_k, \phi_k(i, j) \rangle = \theta_{k,i,j}.$$

We also introduce the shorthands

$$\begin{aligned} a_{1,j} &:= \log \alpha_1(j) = \theta_{1,1,j} \\ a_{k,j} &:= \log \alpha_k(j) = \max_{i \in [M]} \theta_{k,i,j} + a_{k-1,i} \end{aligned}$$

and

$$q_{k,j} := \operatorname{argmax}_{i \in [M]} \theta_{k,i,j} + a_{k-1,i}.$$

Our goal is to compute the gradient w.r.t. $\boldsymbol{\theta} \in \mathbb{R}^{K \times M \times M}$ of

$$\log Z = A = \max_{j \in [M]} a_{K,j}.$$

The soft argmax counterpart of this quantity is

$$Q := \operatorname{argmax}_{j \in [M]} a_{K,j} \in \Delta^M,$$

where we used $\mathcal{P}([M]) = \Delta^M$.

Computing the gradient of A is similar to computing the gradient of a feedforward network, in the sense that $\boldsymbol{\theta}_k$ influences not only a_k but also a_{k+1}, \dots, a_K . Let us introduce the adjoint variable

$$r_{k,i} := \frac{\partial A}{\partial a_{k,i}},$$

which we initialize as

$$r_{K,i} = \frac{\partial A}{\partial a_{K,i}} = Q_i.$$

Since $\theta_{k,i,j}$ directly influences $a_{k,j}$, we have for $k \in [K]$, $i \in [M]$ and $j \in [M]$

$$\begin{aligned}\mu_{k,i,j} &:= \frac{\partial A}{\partial \theta_{k,i,j}} \\ &= \frac{\partial A}{\partial a_{k,j}} \cdot \frac{\partial a_{k,j}}{\partial \theta_{k,i,j}} \\ &= r_{k,j} \cdot q_{k,j,i}.\end{aligned}$$

Since $a_{k,i}$ directly influences $a_{k+1,j}$ for $j \in [M]$, we have for $k \in \{1, \dots, K-1\}$ and $i \in [M]$

$$\begin{aligned}r_{k,i} &= \frac{\partial A}{\partial a_{k,i}} = \sum_{j \in [M]} \frac{\partial A}{\partial a_{k+1,j}} \frac{\partial a_{k+1,j}}{\partial a_{k,i}} \\ &= \sum_{j \in [M]} r_{k+1,j} q_{k+1,j,i} \\ &= \sum_{j \in [M]} \mu_{k+1,i,j}.\end{aligned}$$

We summarize the procedure in Algorithm 10.3. The forward pass uses the softmax operator \max_ε and the softargmax operator $\operatorname{argmax}_\varepsilon$. In the hard max case, in Algorithm 10.2, we used q to store backpointers from integer to integer. In the soft max case, in Algorithm 10.3, we used \mathbf{q} to store **soft backpointers**, that is, discrete probability distributions. In the zero-temperature limit, backpropagation outputs a binary encoding of the solution of backtracking.

Algorithm 10.3 Inference on a chain as backprop with max operators

Input: $\theta \in \mathbb{R}^{K \times M \times M}$

Max operator: \max_{ε}

- 1: Initialize $a_{1,j} := \theta_{1,1,j} \ \forall j \in [M]$
- 2: **for** $k := 2, \dots, K$ **do** \triangleright Forward pass
- 3: **for** $j \in [M]$ **do**
- 4: $a_{k,j} := \max_{\varepsilon} \theta_{k,i,j} + a_{k-1,j} \in \mathbb{R}$
 $\qquad\qquad\qquad i \in [M]$
- 5: $q_{k,j} := \operatorname{argmax}_{\varepsilon} \theta_{k,i,j} + a_{k-1,j} \in \Delta^M$
 $\qquad\qquad\qquad i \in [M]$
- 6: $A := \max_{\varepsilon} a_{K,i} \in \mathbb{R}$
 $\qquad\qquad\qquad i \in [M]$
- 7: $Q := \operatorname{argmax}_{\varepsilon} a_{K,i} \in \Delta^M$
 $\qquad\qquad\qquad i \in [M]$
- 8: Initialize $r_{K,j} = Q_j \ \forall j \in [K]$
- 9: **for** $k := K - 1, \dots, 1$ **do** \triangleright Backward pass
- 10: **for** $i \in [M]$ **do**
- 11: **for** $j \in [M]$ **do**
- 12: $\mu_{k+1,i,j} = r_{k+1,j} \cdot q_{k+1,j,i}$
- 13: $r_{k,i} \leftarrow \mu_{k+1,i,j}$

Outputs: $\max_{\varepsilon} \theta_{1,1,i_1} + \sum_{k=2}^K \theta_{k,i_{k-1},i_k} = A, \nabla A(\theta) = \mu$
 $i_1, \dots, i_K \in [M]^K$

10.10 Summary

- Graphical models represent the conditional dependencies between variables and therefore specify how their joint distribution factorizes.
- There are clear analogies between the worlds of functions and of distributions: the counterparts of computation chains and computation graphs are Markov chains and Bayesian networks.
- Inference on chains and more generally on trees, for exponential family distributions, is equivalent, both statistically and algorithmically, to differentiating the log-partition function.
- The forward-backward algorithm can be seen as using a sum-

product algebra, while the Viterbi algorithm can be seen as using a max-plus algebra. Equivalently, in the log domain, we can see the former as using a soft max, and the latter as using a hard max.

11

Differentiating through optimization

In this chapter, we study how to differentiate through optimization problems, and more generally through nonlinear systems of equations.

11.1 Implicit functions

Implicit functions are functions that do not enjoy an explicit decomposition into elementary functions, for which automatic differentiation, as studied in Chapter 8, can therefore not be directly applied. We describe in this chapter techniques to differentiate through such functions and how to integrate them into an autodiff framework.

Formally, we will denote an implicit function by $\boldsymbol{w}^*(\boldsymbol{\lambda})$, where $\boldsymbol{w}^*: \Lambda \rightarrow \mathcal{W}$. One question is then how to compute the Jacobian $\partial \boldsymbol{w}^*(\boldsymbol{\lambda})$. As a first application one can consider **sensitivity analysis** of a system. For example, $\boldsymbol{w}^*(\boldsymbol{\lambda})$ could correspond to the equilibrium state of a physical system and in this case, $\partial \boldsymbol{w}^*(\boldsymbol{\lambda})$ would tell us about the sensitivity of the system to some parameters $\boldsymbol{\lambda} \in \Lambda$.

11.1.1 Optimization problems

Another example is a function implicitly defined as the solution (assumed unique) of an optimization problem

$$\mathbf{w}^*(\boldsymbol{\lambda}) := \arg \max_{\mathbf{w} \in \mathcal{W}} f(\mathbf{w}, \boldsymbol{\lambda}),$$

where $f: \mathcal{W} \times \Lambda \rightarrow \mathbb{R}$ and \mathcal{W} denotes a constraint set. Note that we use an $\arg \max$ for convenience, but the same applies when using an $\arg \min$.

11.1.2 Nonlinear equations

More generally, $\mathbf{w}^*(\boldsymbol{\lambda})$ can be defined as the root of some function $F: \mathcal{W} \times \Lambda \rightarrow \mathcal{W}$, i.e., $\mathbf{w}^*(\boldsymbol{\lambda})$ is implicitly defined as the function satisfying the (potentially nonlinear) system of equations

$$F(\mathbf{w}, \boldsymbol{\lambda}) = \mathbf{0}$$

for all $\boldsymbol{\lambda} \in \Lambda$.

11.1.3 Application to bilevel optimization

Besides sensitivity analysis, another example of application is **bilevel optimization**. Many times, we want to minimize a function defined as the composition of a fixed function and the solution of an optimization problem. Formally, let $f, g: \mathcal{W} \times \Lambda \rightarrow \mathbb{R}$. We consider the composition $h(\boldsymbol{\lambda})$ defined as

$$h(\boldsymbol{\lambda}) := g(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}), \quad \text{where} \quad \mathbf{w}^*(\boldsymbol{\lambda}) := \arg \max_{\mathbf{w} \in \mathcal{W}} f(\mathbf{w}, \boldsymbol{\lambda}). \quad (11.1)$$

This includes for instance hyperparameter optimization, where f is an inner log-likelihood objective, g is an outer validation loss, $\mathbf{w} \in \mathcal{W}$ are model parameters and $\boldsymbol{\lambda} \in \Lambda$ are model hyperparameters, such as regularization strength, as illustrated in Fig. 11.1. To minimize $h(\boldsymbol{\lambda})$ one generally resorts to a gradient descent scheme w.r.t. $\boldsymbol{\lambda}$, which requires computing $\nabla h(\boldsymbol{\lambda})$. Assuming that $\mathbf{w}^*(\boldsymbol{\lambda})$ is differentiable at $\boldsymbol{\lambda}$, by the chain rule, we obtain the Jacobian

$$\partial h(\boldsymbol{\lambda}) = \partial_1 g(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}) \partial \mathbf{w}^*(\boldsymbol{\lambda}) + \partial_2 g(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}).$$

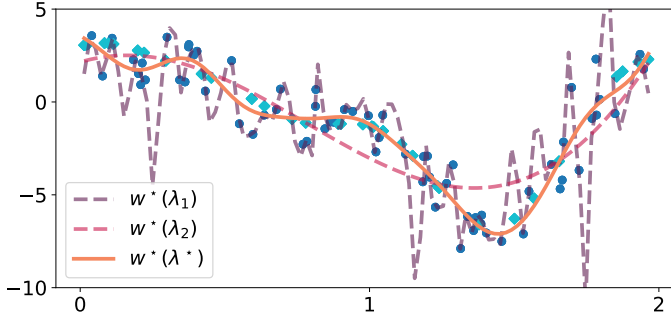


Figure 11.1: Hyperparameter optimization in nonlinear regression can be cast as a bi-level optimization problem. Each line corresponds to the estimator obtained by fitting some training data (in blue circles) using a different hyperparameter λ . Formally, denoting f the training objective, the estimators are $w^*(\lambda) := \arg \min_w f(w; \lambda)$. The goal is to find the best hyperparameter that fits some validation data (here in cyan diamonds), that is, minimizing $h(\lambda) := g(w^*(\lambda), \lambda)$, where g is the validation objective. A too small λ_1 leads to overfitting the training objective and performs badly on validation objective. Conversely, a larger λ_2 underfits both training and validation objectives. The optimal parameter λ^* minimizes the validation objective and may be obtained by iterating gradient descent w.r.t. λ . This requires gradients of $h(\lambda) = g(w^*(\lambda), \lambda)$ w.r.t. λ .

Using $\partial h(\lambda)^\top = \nabla h(\lambda)$ (see Remark 2.4), we obtain the gradient

$$\nabla h(\lambda) = \partial w^*(\lambda)^\top \nabla_1 g(w^*(\lambda), \lambda) + \nabla_2 g(w^*(\lambda), \lambda).$$

The only problematic term is $\partial w^*(\lambda)$, as it requires **argmax** differentiation. Indeed, most of the time, there is no explicit formula for $w^*(\lambda)$ and it does not decompose into elementary functions.

11.2 Envelope theorems

In the special case $g = f$, the composition h defined in Eq. (11.1) is simply given by

$$h(\lambda) = f(w^*(\lambda), \lambda) = \max_{w \in \mathcal{W}} f(w, \lambda).$$

That is, we no longer need **argmax** differentiation, but only **max** differentiation, which, as we shall now see is much easier. The function h is often called a **value function** (Fleming and Rishel, 2012). The

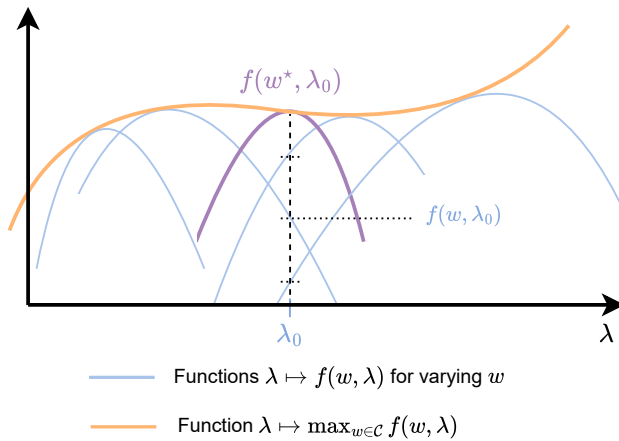


Figure 11.2: The graph of $h(\lambda) = \max_{w \in \mathcal{W}} f(w, \lambda)$ is the upper-envelope of the graphs of the functions $\lambda \mapsto f(w, \lambda)$ for all $w \in \mathcal{W}$.

reason for the name “envelope” is illustrated in Fig. 11.2. We emphasize that there is not one, but several envelope theorems, depending on the assumptions on f .

11.2.1 Danskin’s theorem

When f is concave-convex, we can use Danskin’s theorem.

Theorem 11.1 (Danskin’s theorem). Let $f: \mathcal{W} \times \Lambda \rightarrow \mathbb{R}$ and \mathcal{W} be a compact convex set. Let

$$h(\lambda) := \max_{w \in \mathcal{W}} f(w, \lambda)$$

and

$$w^*(\lambda) := \arg \max_{w \in \mathcal{W}} f(w, \lambda).$$

If f is **concave** in w , **convex** in λ , and the maximum $w^*(\lambda)$ is unique, then the function h is differentiable with gradient

$$\nabla h(\lambda) = \nabla_2 f(w^*(\lambda), \lambda).$$

If the maximum is not unique, we get a subgradient.

Informally, Danskin's theorem means that we can treat $\mathbf{w}^*(\boldsymbol{\lambda})$ as if it were a constant of $\boldsymbol{\lambda}$, i.e., we do not need to differentiate through it, even though it depends on $\boldsymbol{\lambda}$. Danskin's theorem can also be used to differentiate through a minimum, $h(\boldsymbol{\lambda}) = \min_{\mathbf{w} \in \mathcal{W}} f(\mathbf{w}, \boldsymbol{\lambda})$, if $f(\mathbf{w}, \boldsymbol{\lambda})$ is convex in \mathbf{w} and concave in $\boldsymbol{\lambda}$, as we now illustrate.

Example 11.1 (Illustration of Danskin's theorem). Let us define $h(\lambda) := \min_{w \in \mathbb{R}} f(w, \lambda)$, where $f(w, \lambda) := \frac{\lambda}{2}w^2 + bw + c$ and $\lambda > 0$. Let $w^*(\lambda)$ be the minimum. The derivative of f w.r.t. λ is $\frac{1}{2}w^2$. From Danskin's theorem, we have $h'(\lambda) = \frac{1}{2}w^*(\lambda)$. Let us check that this result is correct. The derivative of f w.r.t. w is $\lambda w + b$. Setting it to zero, we get $w^*(\lambda) = -\frac{b}{\lambda}$. We thus obtain $h'(\lambda) = \frac{1}{2}\frac{b^2}{\lambda^2}$. Plugging $w^*(\lambda)$ back into $f(w, \lambda)$, we get $h(\lambda) = -\frac{1}{2}\frac{b^2}{\lambda} + c$. Using $(\frac{1}{\lambda})' = -\frac{1}{\lambda^2}$, we indeed obtain the same result for $h'(\lambda)$.

Danskin's theorem has a simple interpretation for functions that are linear in $\boldsymbol{\lambda}$ as shown below.

Example 11.2 (Convex conjugate). Let $f(\mathbf{w}, \boldsymbol{\lambda}) := \langle \mathbf{w}, \boldsymbol{\lambda} \rangle - \Omega(\mathbf{w})$ with Ω convex. We then have $h(\boldsymbol{\lambda}) = \max_{\mathbf{w} \in \mathcal{W}} \langle \mathbf{w}, \boldsymbol{\lambda} \rangle - \Omega(\mathbf{w}) =: \Omega^*(\boldsymbol{\lambda})$, where Ω^* denotes the convex conjugate of Ω . Since f satisfies the conditions of Danskin's theorem and since we have $\nabla_{\mathbf{w}} f(\mathbf{w}, \boldsymbol{\lambda}) = \mathbf{w}$, we obtain $\nabla h(\boldsymbol{\lambda}) = \nabla \Omega^*(\boldsymbol{\lambda}) = \mathbf{w}^*(\boldsymbol{\lambda})$. In other words, in this special case, the gradient of the max is equal to the argmax. This is due to the fact that $f(\mathbf{w}, \boldsymbol{\lambda})$ is linear in $\boldsymbol{\lambda}$.

Another application is saddle point optimization.

Example 11.3 (Saddle point problem). Consider the saddle point problem $\min_{\boldsymbol{\lambda} \in \Lambda} \max_{\mathbf{w} \in \mathcal{W}} f(\mathbf{w}, \boldsymbol{\lambda})$. If it is difficult to minimize w.r.t. $\boldsymbol{\lambda}$ but easy to maximize w.r.t. \mathbf{w} , we can rewrite the problem as $\min_{\boldsymbol{\lambda} \in \Lambda} h(\boldsymbol{\lambda})$, where $h(\boldsymbol{\lambda}) := \max_{\mathbf{w} \in \mathcal{W}} f(\mathbf{w}, \boldsymbol{\lambda})$, and use $\nabla h(\boldsymbol{\lambda})$ to perform (projected) gradient descent w.r.t. $\boldsymbol{\lambda}$.

11.2.2 Rockafellar's theorem

A related theorem can be proved under different assumptions on f , in particular without concavity w.r.t. \mathbf{w} .

Theorem 11.2 (Rockafellar's envelope theorem). Let $f: \mathcal{W} \times \Lambda \rightarrow \mathbb{R}$ and \mathcal{W} be a compact convex set. Let

$$h(\boldsymbol{\lambda}) := \max_{\boldsymbol{w} \in \mathcal{W}} f(\boldsymbol{w}, \boldsymbol{\lambda})$$

and

$$\boldsymbol{w}^*(\boldsymbol{\lambda}) := \arg \max_{\boldsymbol{w} \in \mathcal{W}} f(\boldsymbol{w}, \boldsymbol{\lambda}).$$

If f is continuously differentiable in $\boldsymbol{\lambda}$ for all $\boldsymbol{w} \in \mathcal{W}$, $\nabla_1 f$ is continuous and the maximum $\boldsymbol{w}^*(\boldsymbol{\lambda})$ is unique, then the function h is differentiable with gradient

$$\nabla h(\boldsymbol{\lambda}) = \nabla_2 f(\boldsymbol{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}).$$

See Rockafellar and Wets (2009, Theorem 10.31). Compared to Danskin's theorem, Rockafellar's theorem does not require f to be concave-convex, but requires stronger assumptions on the differentiability of f .

11.3 Implicit function theorem

11.3.1 Univariate functions

The implicit function theorem (IFT) provides conditions under which an implicit relationship of the form $F(w, \lambda) = 0$ can be rewritten as a function $w = w^*(\lambda)$ locally, and provides a way to compute its derivative w.r.t. λ .

Theorem 11.3 (Implicit function theorem, univariate case). Let $F: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$. Assume $F(w, \lambda)$ is a continuously differentiable function in a neighborhood \mathcal{U} of (w_0, λ_0) such that $F(w_0, \lambda_0) = 0$ and $\partial_1 F(w_0, \lambda_0) \neq 0$. Then there exists a neighborhood $\mathcal{V} \subseteq \mathcal{U}$ of (w_0, λ_0) in which there is a function $w^*(\lambda)$ such that

- $w^*(\lambda_0) = w_0$,
- $F(w^*(\lambda), \lambda) = 0$ for all λ in the neighborhood \mathcal{V} ,

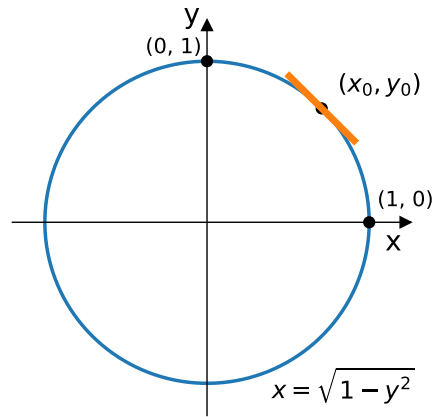


Figure 11.3: The circle equation $F(x, y) := x^2 + y^2 - 1 = 0$ is not a function from $y \in [-1, 1]$ to $x \in [-1, 1]$, as there are always two possible x values, $x = \sqrt{1 - y^2}$ or $x = -\sqrt{1 - y^2}$. However, locally around some point (x_0, y_0) , e.g., such that $x_0 > 0$ and $y_0 > 0$ (upper-right quadrant), the function $x = x^*(y) = \sqrt{1 - y^2}$ is well-defined. The implicit function theorem gives conditions for such a function to exist locally and provides a way to compute its derivative.

$$\bullet \quad \partial w^*(\lambda) = -\frac{\partial_2 F(w^*(\lambda), \lambda)}{\partial_1 F(w^*(\lambda), \lambda)}.$$

We postpone the proof to the multivariate case and begin with a classical example of application of the theorem.

Example 11.4 (Equation of the unit circle). We use $w \equiv x$ and $\lambda \equiv y$ for clarity. Let $F(x, y) := x^2 + y^2 - 1$. In general, we cannot rewrite the unit circle equation $F(x, y) = 0$ as a function from y to x , because for every $y \in [-1, 1]$, there are always two possible x values, namely, $x = \sqrt{1 - y^2}$ or $x = -\sqrt{1 - y^2}$. However, locally around some point (x_0, y_0) , e.g., such that $x_0 > 0$ and $y_0 > 0$ (upper-right quadrant), the function $x = x^*(y) = \sqrt{1 - y^2}$ is well-defined. Using $\partial_1 F(x, y) = 2x$ and $\partial_2 F(x, y) = 2y$, we get $\partial x^*(y) = -\frac{\partial_2 F(x^*(y), y)}{\partial_1 F(x^*(y), y)} = -\frac{2y}{2x^*(y)} = -\frac{y}{\sqrt{1 - y^2}}$ in that neighborhood (the upper right quadrant in this case). This is indeed the same derivative expression as if we used the chain rule on $\sqrt{1 - y^2}$ and is well-defined on $y \in [0, 1)$.

In the above simple example, we can easily derive an explicit function relating y to x in a given neighborhood, but this is not always the case. The IFT gives us conditions guaranteeing that such function **exists** and a way to **differentiate** it, but not a way to **construct** such a function. In fact, finding $w^*(\lambda)$ such that $F(w^*(\lambda), \lambda) = 0$ typically involves a root finding algorithm, an optimization algorithm, a nonlinear system solver, etc.

Example 11.5 (Polynomial). Let $F(w, \lambda) = w^5 + w^3 + w - \lambda$. According to the Abel-Ruffini theorem (Tignol, 2015), quintics (polynomials of degree 5) do not enjoy roots in terms of radicals and one must resort to numerical root finding. In addition, odd-degree polynomials have real roots. Moreover, $\partial_1 F(w, \lambda) = 5w^4 + 3w^2 + 1$ is strictly positive. Therefore, by the intermediate value theorem, there must be only one root $w^*(\lambda)$ such that $F(w^*(\lambda), \lambda) = 0$. This unique root can for example be found by bisection. Using the IFT, its derivative is found to be $\partial w^*(\lambda) = (5w^*(\lambda)^4 + 3w^*(\lambda)^2 + 1)^{-1}$.

While an implicit function is differentiable at a point if the assumptions of the IFT hold in a neighborhood of that point, the reciprocal is not true: failure of the IFT assumptions does not necessarily mean that the implicit function is not differentiable, as we now illustrate.

Example 11.6 (IFT conditions are not necessary for differentiability). Consider $F(w, \lambda) = (w - \lambda)^2$. We clearly have that $F(w^*(\lambda), \lambda) = 0$ if we define $w^*(\lambda) = \lambda$, the identity function. It is clearly differentiable for all λ , yet the assumptions of the IFT fail, since we have $\partial_1 F(w, \lambda) = 2(w - \lambda)$ and therefore $\partial_1 F(0, 0) = 0$.

11.3.2 Multivariate functions

We now present the IFT in the general multivariate setting. Informally, if $F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}) = \mathbf{0}$, then by the chain rule, we have

$$\partial_1 F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}) \partial \mathbf{w}^*(\boldsymbol{\lambda}) + \partial_2 F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}) = \mathbf{0},$$

meaning that the Jacobian $\partial \mathbf{w}^*(\boldsymbol{\lambda})$, assuming that it exists, satisfies

$$-\partial_1 F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}) \partial \mathbf{w}^*(\boldsymbol{\lambda}) = \partial_2 F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}).$$

The IFT gives us conditions for the existence of $\partial \mathbf{w}^*(\boldsymbol{\lambda})$.

Theorem 11.4 (Implicit function theorem, multivariate case). Let us define $F: \mathcal{W} \times \Lambda \rightarrow \mathcal{W}$. Assume $F(\mathbf{w}, \boldsymbol{\lambda})$ is a continuously differentiable function in a neighborhood of $(\mathbf{w}_0, \boldsymbol{\lambda}_0)$ such that $F(\mathbf{w}_0, \boldsymbol{\lambda}_0) = \mathbf{0}$ and $\partial_1 F(\mathbf{w}_0, \boldsymbol{\lambda}_0)$ is invertible, i.e., its determinant is nonzero. Then there exists a neighborhood of $\boldsymbol{\lambda}_0$ in which there is a function $\mathbf{w}^*(\boldsymbol{\lambda})$ such that

- $\mathbf{w}^*(\boldsymbol{\lambda}_0) = \mathbf{w}_0$,
- $F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}) = \mathbf{0}$ for all $\boldsymbol{\lambda}$ in the neighborhood,
- $-\partial_1 F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}) \partial \mathbf{w}^*(\boldsymbol{\lambda}) = \partial_2 F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda})$
 $\iff \partial \mathbf{w}^*(\boldsymbol{\lambda}) = -\partial_1 F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda})^{-1} \partial_2 F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda})$.

We begin with a simple unconstrained optimization algorithm.

Example 11.7 (Unconstrained optimization). Assume we want to differentiate through $\mathbf{w}^*(\boldsymbol{\lambda}) = \arg \min_{\mathbf{w} \in \mathbb{R}^P} f(\mathbf{w}, \boldsymbol{\lambda})$, where f is strictly convex in \mathbf{w} , which ensures that the solution is unique. From the stationary conditions, if we define $F(\mathbf{w}, \boldsymbol{\lambda}) := \nabla_1 f(\mathbf{w}, \boldsymbol{\lambda})$, then $\mathbf{w}^*(\boldsymbol{\lambda})$ is uniquely characterized as the root of F in the first argument, i.e., $F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}) = \mathbf{0}$. We have $\partial_1 F(\mathbf{w}, \boldsymbol{\lambda}) = \nabla_1^2 f(\mathbf{w}, \boldsymbol{\lambda})$, the Hessian of f in \mathbf{w} , and $\partial_2 F(\mathbf{w}, \boldsymbol{\lambda}) = \partial_2 \nabla_1 f(\mathbf{w}, \boldsymbol{\lambda})$, the cross derivatives of f in \mathbf{w} and $\boldsymbol{\lambda}$. Therefore, assuming that the Hessian is well-defined and invertible at $(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda})$, we can use the IFT to differentiate through $\mathbf{w}^*(\boldsymbol{\lambda})$ and obtain $\partial \mathbf{w}^*(\boldsymbol{\lambda}) = -(\nabla_1^2 f(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}))^{-1} \partial_2 \nabla_1 f(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda})$.

Next, we generalize the previous example, by allowing constraints in the optimization problem.

Example 11.8 (Constrained optimization). Now, assume we want to differentiate through $\mathbf{w}^*(\boldsymbol{\lambda}) = \arg \min_{\mathbf{w} \in \mathcal{C}} f(\mathbf{w}, \boldsymbol{\lambda})$, where f is strictly convex in \mathbf{w} and $\mathcal{C} \subseteq \mathcal{W}$ is a convex set. A solution is characterized by the fixed point equation $\mathbf{w}^*(\boldsymbol{\lambda}) = P_{\mathcal{C}}(\mathbf{w}^*(\boldsymbol{\lambda}) -$

$\eta \nabla_1 f(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}))$, for any $\eta > 0$, where $P(\mathbf{y}) := \arg \min_{\mathbf{x} \in \mathcal{C}} \|\mathbf{x} - \mathbf{y}\|_2^2$ is the Euclidean projection of \mathbf{y} onto \mathcal{C} . Therefore, $\mathbf{w}^*(\boldsymbol{\lambda})$ is the root of $F(\mathbf{w}, \boldsymbol{\lambda}) = \mathbf{w} - P_{\mathcal{C}}(\mathbf{w} - \eta \nabla_1 f(\mathbf{w}, \boldsymbol{\lambda}))$ (see Chapter 16). We can differentiate through $\mathbf{w}^*(\boldsymbol{\lambda})$ using the IFT, assuming that the conditions of the theorem apply. Note that $\partial_1 F(\mathbf{w}, \boldsymbol{\lambda})$ requires the expression of the Jacobian $\partial P_{\mathcal{C}}(\mathbf{y})$. Fortunately, $P_{\mathcal{C}}(\mathbf{y})$ and its Jacobian are easy to compute for many sets \mathcal{C} (Blondel *et al.*, 2021).

11.3.3 JVP and VJP of implicit functions

To integrate an implicit function $\mathbf{w}^*(\boldsymbol{\lambda})$ in an autodiff framework, we need to be able to compute its JVP or VJP. This is the purpose of the next proposition.

Proposition 11.1 (JVP and VJP of implicit functions). Let $\mathbf{w}^*: \Lambda \rightarrow \mathcal{W}$ be a function implicitly defined as the solution of $F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}) = \mathbf{0}$, for some function $F: \mathcal{W} \times \Lambda \rightarrow \mathcal{W}$. Define

$$\begin{aligned} A &:= -\partial_1 F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}) \\ B &:= \partial_2 F(\mathbf{w}^*(\boldsymbol{\lambda}), \boldsymbol{\lambda}). \end{aligned}$$

Assume the assumptions of the IFT hold. The JVP $\mathbf{t} := \partial \mathbf{w}^*(\boldsymbol{\lambda}) \mathbf{v}$ in the input direction $\mathbf{v} \in \Lambda$ is obtained by solving the linear system

$$A\mathbf{t} = B\mathbf{v}.$$

The VJP $\partial \mathbf{w}^*(\boldsymbol{\lambda})^* \mathbf{u}$ in the output direction $\mathbf{u} \in \mathcal{W}$ is obtained by solving the linear system

$$A^* \mathbf{r} = \mathbf{u}.$$

Using the solution \mathbf{r} , we get

$$\partial \mathbf{w}^*(\boldsymbol{\lambda})^* \mathbf{u} = \partial \mathbf{w}^*(\boldsymbol{\lambda})^* A^* \mathbf{r} = B^* \mathbf{r}.$$

Note that in the above linear systems, we can access to A and B as linear maps, the JVPs of F . Their adjoints, A^* and B^* , correspond to the VJPs of F . To solve these systems, we can therefore use **matrix-free** solvers as detailed in Section 9.4. For example, when A is symmetric pos-

itive semi-definite, we can use the conjugate gradient method (Hestenes, Stiefel, *et al.*, 1952). When A is not symmetric positive definite, we can use GMRES (Saad and Schultz, 1986) or BiCGSTAB (Vorst and Vorst, 1992).

11.3.4 Proof of the implicit function theorem

We prove the theorem using the inverse function theorem presented in Theorem 11.5. Define

$$f(\boldsymbol{\lambda}, \mathbf{w}) = (\boldsymbol{\lambda}, F(\mathbf{w}, \boldsymbol{\lambda}))$$

which goes from $\mathbb{R}^Q \times \mathbb{R}^P$ onto $\mathbb{R}^Q \times \mathbb{R}^P$. The Jacobian of f is

$$\partial f(\boldsymbol{\lambda}, \mathbf{w}) = \begin{pmatrix} \mathbf{I} & 0 \\ \partial_2 F(\mathbf{w}, \boldsymbol{\lambda}) & \partial_1 F(\mathbf{w}, \boldsymbol{\lambda}) \end{pmatrix}.$$

So at $\mathbf{w}_0, \boldsymbol{\lambda}_0$, we have $\det(\partial f(\boldsymbol{\lambda}_0, \mathbf{w}_0)) = \det(\mathbf{I}) \det(\partial_1 F(\mathbf{w}_0, \boldsymbol{\lambda}_0)) > 0$ since we assumed $\partial_1 F(\mathbf{w}_0, \boldsymbol{\lambda}_0)$ invertible. By the inverse function theorem, the function f is then invertible in a neighborhood N of $f(\boldsymbol{\lambda}_0, \mathbf{w}_0) = (\boldsymbol{\lambda}_0, \mathbf{0})$. In particular, it is invertible in $N \cap \{(\boldsymbol{\lambda}, \mathbf{0}), \boldsymbol{\lambda} \in \mathbb{R}^Q\}$. The solution of the implicit equation in a neighborhood of $\boldsymbol{\lambda}_0$ is then $(\boldsymbol{\lambda}, \mathbf{w}^*(\boldsymbol{\lambda})) = f^{-1}(\boldsymbol{\lambda}, \mathbf{0})$. By the inverse function theorem, f^{-1} is continuously differentiable inverse and so is $\mathbf{w}^*(\boldsymbol{\lambda})$. The derivative $\partial \mathbf{w}^*(\boldsymbol{\lambda})$ from the differential of the inverse as

$$\begin{pmatrix} \sim & \sim \\ \partial \mathbf{w}^*(\boldsymbol{\lambda}) & \sim \end{pmatrix} = \partial f^{-1}(\boldsymbol{\lambda}, \mathbf{0}),$$

and by the inverse function theorem, we have $\partial f^{-1}(\boldsymbol{\lambda}, \mathbf{0}) = (\partial f(\boldsymbol{\lambda}, \mathbf{w}^*(\boldsymbol{\lambda})))^{-1}$. So using block matrix inversions formula

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \sim & \sim \\ -(D - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1} & \sim \end{pmatrix},$$

we get the claimed expression. Though we expressed the proof in terms of Jacobians and matrices, the result naturally holds for the corresponding linear operators, JVPs, VJPs, and their inverses.

11.4 Adjoint state method

11.4.1 Differentiating nonlinear equations

We describe in this section the adjoint state method (a.k.a. adjoint method, method of adjoints, adjoint sensitivity method). The method can be used to compute the gradient of the composition of an explicit function and an **implicit function**, defined through an **equality constraint** (e.g., a **nonlinear equation**). The method dates back to C  a (1986).

Suppose a variable $\mathbf{s} \in \mathcal{S}$ (which corresponds to a **state** in optimal control) is implicitly defined given some parameters $\mathbf{w} \in \mathcal{W}$ through the (potentially nonlinear) equation $c(\mathbf{s}, \mathbf{w}) = \mathbf{0}$, where $c: \mathcal{S} \times \mathcal{W} \rightarrow \mathcal{S}$. Assuming \mathbf{s} is **uniquely** determined for all $\mathbf{w} \in \mathcal{W}$, this defines an **implicit function** $\mathbf{s}^*(\mathbf{w})$ from \mathcal{W} to \mathcal{S} such that $c(\mathbf{s}^*(\mathbf{w}), \mathbf{w}) = \mathbf{0}$. Given an objective function $L: \mathcal{S} \times \mathcal{W} \rightarrow \mathbb{R}$, the goal of the adjoint state method is then to compute the gradient of

$$L(\mathbf{w}) := L(\mathbf{s}^*(\mathbf{w}), \mathbf{w}).$$

However, this is not trivial as $\mathbf{s}^*(\mathbf{w})$ is an implicit function. For instance, this can be used to convert the **equality-constrained** problem

$$\min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{s}, \mathbf{w}) \quad \text{s.t.} \quad c(\mathbf{s}, \mathbf{w}) = \mathbf{0}.$$

into the **unconstrained** problem

$$\min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{s}^*(\mathbf{w}), \mathbf{w}).$$

Access to $\nabla L(\mathbf{w})$ allows us to solve this problem by gradient descent.

Proposition 11.2 (Adjoint state method). Let $c: \mathcal{S} \times \mathcal{W} \rightarrow \mathcal{S}$ be a mapping defining constraints of the form $c(\mathbf{s}, \mathbf{w})$. Assume that for each $\mathbf{w} \in \mathcal{W}$, there exists a unique $\mathbf{s}^*(\mathbf{w})$ satisfying $c(\mathbf{s}^*(\mathbf{w}), \mathbf{w}) = \mathbf{0}$ and that $\mathbf{s}^*(\mathbf{w})$ is differentiable. The gradient of

$$L(\mathbf{w}) := L(\mathbf{s}^*(\mathbf{w}), \mathbf{w}),$$

for some differentiable function $L: \mathcal{S} \times \mathcal{W} \rightarrow \mathbb{R}$, is given by

$$\nabla L(\mathbf{w}) = \nabla_2 L(\mathbf{s}^*(\mathbf{w}), \mathbf{w}) + \partial_2 c(\mathbf{s}^*(\mathbf{w}), \mathbf{w})^* \mathbf{r}^*(\mathbf{w}),$$

where $\mathbf{r}^*(\mathbf{w})$ is the solution of the linear system

$$\partial_1 c(\mathbf{s}^*(\mathbf{w}), \mathbf{w})^* \mathbf{r} = -\nabla_1 L(\mathbf{s}^*(\mathbf{w}), \mathbf{w}).$$

As shown in the proof below, $\mathbf{r}^*(\mathbf{w})$ corresponds to a **Lagrange multiplier**. The linear system can be solved using matrix-free solvers.

11.4.2 Relation with envelope theorems

Because \mathbf{s} is uniquely determined for any $\mathbf{w} \in \mathcal{W}$ by $c(\mathbf{s}, \mathbf{w}) = \mathbf{0}$, we can alternatively rewrite $L(\mathbf{w})$ as the trivial minimization or maximization,

$$\begin{aligned} L(\mathbf{w}) &= \min_{\mathbf{s} \in \mathcal{S}} L(\mathbf{s}, \mathbf{w}) \quad \text{s.t.} \quad c(\mathbf{s}, \mathbf{w}) = \mathbf{0} \\ &= \max_{\mathbf{s} \in \mathcal{S}} L(\mathbf{s}, \mathbf{w}) \quad \text{s.t.} \quad c(\mathbf{s}, \mathbf{w}) = \mathbf{0}. \end{aligned}$$

Therefore, the adjoint state method can be seen as an envelope theorem for computing $\nabla L(\mathbf{w})$, for the case when \mathbf{w} is involved in **both** the objective function and in the **equality constraint**.

11.4.3 Proof using the method of Lagrange multipliers

Classically, the adjoint state method is derived using the method of Lagrange multipliers. Let us introduce the **Lagrangian** associated with L and c ,

$$\mathcal{L}(\mathbf{s}, \mathbf{w}, \mathbf{r}) := L(\mathbf{s}, \mathbf{w}) + \langle \mathbf{r}, c(\mathbf{s}, \mathbf{w}) \rangle,$$

where $\mathbf{r} \in \mathcal{S}$ is the **Lagrange multiplier** associated with the equality constraint $c(\mathbf{s}, \mathbf{w}) = \mathbf{0}$. In the optimal control literature, \mathbf{r} is often called the **adjoint variable** or **adjoint state**. The gradients of the Lagrangian are

$$\begin{aligned} \nabla_{\mathbf{s}} \mathcal{L}(\mathbf{s}, \mathbf{w}, \mathbf{r}) &= \nabla_1 L(\mathbf{s}, \mathbf{w}) + \partial_1 c(\mathbf{s}, \mathbf{w})^* \mathbf{r} \\ \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{s}, \mathbf{w}, \mathbf{r}) &= \nabla_2 L(\mathbf{s}, \mathbf{w}) + \partial_2 c(\mathbf{s}, \mathbf{w})^* \mathbf{r} \\ \nabla_{\mathbf{r}} \mathcal{L}(\mathbf{s}, \mathbf{w}, \mathbf{r}) &= c(\mathbf{s}, \mathbf{w}), \end{aligned}$$

where $\partial_i c(\mathbf{s}, \mathbf{w})^*$ are the **adjoint operators**. Setting $\nabla_{\mathbf{r}} \mathcal{L}(\mathbf{s}, \mathbf{w}, \mathbf{r})$ to zero gives the constraint $c(\mathbf{s}, \mathbf{w}) = \mathbf{0}$. Setting $\nabla_{\mathbf{s}} \mathcal{L}(\mathbf{s}, \mathbf{w}, \mathbf{r})$ to zero gives the so-called **adjoint state equation**

$$\partial_1 c(\mathbf{s}, \mathbf{w})^* \mathbf{r} = -\nabla_1 L(\mathbf{s}, \mathbf{w}).$$

Solving this linear system w.r.t. \mathbf{r} at $\mathbf{s} = \mathbf{s}^*(\mathbf{w})$ gives the adjoint variable $\mathbf{r}^*(\mathbf{w})$. We then get

$$\begin{aligned}\nabla L(\mathbf{w}) &= \nabla_2 \mathcal{L}(\mathbf{s}^*(\mathbf{w}), \mathbf{w}, \mathbf{r}^*(\mathbf{w})) \\ &= \nabla_2 L(\mathbf{s}^*(\mathbf{w}), \mathbf{w}) + \partial_2 c(\mathbf{s}^*(\mathbf{w}), \mathbf{w})^* \mathbf{r}^*(\mathbf{w}),\end{aligned}$$

which concludes the proof.

11.4.4 Proof using the implicit function theorem

A more direct proof is possible thanks to the implicit function theorem (Section 11.3). Using the chain rule, we get

$$\nabla L(\mathbf{w}) = \nabla_2 L(\mathbf{s}^*(\mathbf{w}), \mathbf{w}) + \partial \mathbf{s}^*(\mathbf{w})^* \nabla_1 L(\mathbf{s}^*(\mathbf{w}), \mathbf{w}),$$

where $\partial \mathbf{s}^*(\mathbf{w})^*$ is the VJP of \mathbf{s}^* , a linear map from \mathcal{S} to \mathcal{W} .

Computationally, the main difficulty is to apply $\partial \mathbf{s}^*(\mathbf{w})^*$ to the vector $\mathbf{u} = \nabla_1 L(\mathbf{s}^*(\mathbf{w}), \mathbf{w}) \in \mathcal{S}$. Using the implicit function theorem (Section 11.3) on the implicit function $c(\mathbf{s}^*(\mathbf{w}), \mathbf{w}) = \mathbf{0}$, and Proposition 11.1, we get the linear system $A^* \mathbf{r} = \mathbf{u}$, where $A^* := \partial_1 c(\mathbf{s}^*(\mathbf{w}), \mathbf{w})^*$ is a linear map from \mathcal{S} to \mathcal{S} . After solving for \mathbf{r} , we get $\partial \mathbf{s}^*(\mathbf{w})^* \mathbf{u} = B^* \mathbf{r}$, where $B^* := \partial_2 c(\mathbf{s}^*(\mathbf{w}), \mathbf{w})^*$ is a linear map from \mathcal{S} to \mathcal{W} . Putting everything together, we get

$$\nabla L(\mathbf{w}) = \nabla_2 L(\mathbf{s}^*(\mathbf{w}), \mathbf{w}) + \partial_2 c(\mathbf{s}^*(\mathbf{w}), \mathbf{w})^* \mathbf{r}.$$

11.4.5 Reverse mode as adjoint method with backsubstitution

In this section, we revisit reverse-mode autodiff from the perspective of the adjoint state method. For clarity, we focus our exposition on feedforward networks with input $\mathbf{x} \in \mathcal{X}$ and network weights $\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_K) \in \mathcal{W}_1 \times \dots \times \mathcal{W}_K$,

$$\begin{aligned}\mathbf{s}_0 &:= \mathbf{x} \in \mathcal{X} \\ \mathbf{s}_1 &:= f_1(\mathbf{s}_0, \mathbf{w}_1) \in \mathcal{S}_1 \\ &\vdots \\ \mathbf{s}_K &:= f_K(\mathbf{s}_{K-1}, \mathbf{w}_K) \in \mathcal{S}_K \\ f(\mathbf{w}) &:= \mathbf{s}_K.\end{aligned}\tag{11.2}$$

Here we focus on gradients with respect to the parameters \mathbf{w} , hence the notation $f(\mathbf{w})$. We can use the adjoint state method to recover reverse-mode autodiff, and prove its **correctness** in the process. While we focus for simplicity on feedforward networks, our exposition can be generalized to computation graphs.

Feedforward networks as the solution of a nonlinear equation

While we defined the set of intermediate computations $\mathbf{s} = (\mathbf{s}_1, \dots, \mathbf{s}_K) \in \mathcal{S}_1 \times \dots \times \mathcal{S}_K$ as a sequence of operations, they can also be defined as the unique solution of the **nonlinear equation** $c(\mathbf{s}, \mathbf{w}) = \mathbf{0}$, where

$$c(\mathbf{s}, \mathbf{w}) := \begin{pmatrix} \mathbf{s}_1 - f_1(\mathbf{x}, \mathbf{w}_1) \\ \mathbf{s}_2 - f_2(\mathbf{s}_1, \mathbf{w}_2) \\ \vdots \\ \mathbf{s}_K - f_K(\mathbf{s}_{K-1}, \mathbf{w}_K) \end{pmatrix}.$$

This defines an **implicit function** $\mathbf{s}^*(\mathbf{w}) = (\mathbf{s}_1^*(\mathbf{w}), \dots, \mathbf{s}_K^*(\mathbf{w}))$, the solution of this nonlinear system, which is given by the variables $\mathbf{s}_1, \dots, \mathbf{s}_K$ defined in Eq. (11.2). The output of the feedforward network is then $f(\mathbf{w}) = \mathbf{s}_K^*(\mathbf{w})$.

In machine learning, the final layer $\mathbf{s}_K^*(\mathbf{w})$ is typically fed into a loss ℓ , to define

$$L(\mathbf{w}) := \ell(\mathbf{s}_K^*(\mathbf{w}); \mathbf{y}).$$

Note that an alternative is to write $L(\mathbf{w})$ as

$$L(\mathbf{w}) = \min_{\mathbf{s} \in \mathcal{S}} \ell(\mathbf{s}; \mathbf{y}) \quad \text{s.t.} \quad c(\mathbf{s}, \mathbf{w}) = \mathbf{0}.$$

More generally, if we just want to compute the VJP of $\mathbf{s}_K^*(\mathbf{w})$ in some direction $\mathbf{u}_K \in \mathcal{S}_K$, we can define the scalar-valued function

$$L(\mathbf{w}) := \ell(\mathbf{s}_K^*(\mathbf{w}); \mathbf{u}_K) := \langle \mathbf{s}_K^*(\mathbf{w}), \mathbf{u}_K \rangle$$

so that

$$\partial f(\mathbf{w})^* \mathbf{u}_K = \nabla L(\mathbf{w}).$$

Let us define $\mathbf{u} \in \mathcal{S}_1 \times \dots \times \mathcal{S}_{K-1} \times \mathcal{S}_K$ as $\mathbf{u} := (\mathbf{0}, \dots, \mathbf{0}, \nabla_1 \ell(f(\mathbf{w}); \mathbf{y}))$ (gradient of the loss ℓ case) or $\mathbf{u} := (\mathbf{0}, \dots, \mathbf{0}, \mathbf{u}_K)$ (VJP of f in the

direction \mathbf{u}_K case). Using the adjoint state method, we know that the gradient of this objective is obtained as

$$\nabla L(\mathbf{w}) = \partial_2 c(\mathbf{s}(\mathbf{w}), \mathbf{w})^* \mathbf{r}^*(\mathbf{w}),$$

for $\mathbf{r}^*(\mathbf{w})$ the solution of the linear system

$$\partial_1 c(\mathbf{s}(\mathbf{w}), \mathbf{w})^* \mathbf{r} = -\mathbf{u}.$$

Solving the linear system using backsubstitution

The JVP of the constraint function c at $\mathbf{s}^*(\mathbf{w})$, materialized as a matrix, takes the form of a **block lower-triangular matrix**

$$\partial_1 c(\mathbf{s}^*(\mathbf{w}), \mathbf{w}) = \begin{pmatrix} \mathbf{I} & 0 & \dots & \dots & 0 \\ -A_1 & \mathbf{I} & \ddots & & \vdots \\ 0 & -A_2 & \mathbf{I} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & -A_K & \mathbf{I} \end{pmatrix},$$

where $A_k := \partial_1 f_k(\mathbf{s}_{k-1}, \mathbf{w}_k)$. Crucially the triangular structure of the JVP stems from the fact that each intermediate activation only depends from the past intermediate activations. Therefore, the constraints, corresponding to the lines of the Jacobian, cannot introduce non-zero values beyond its diagonal. The VJP takes the form of a **block upper-triangular matrix**

$$\partial_1 c(\mathbf{s}^*(\mathbf{w}), \mathbf{w})^* = \begin{pmatrix} \mathbf{I} & -A_1^* & 0 & \dots & 0 \\ 0 & \mathbf{I} & -A_2^* & \ddots & \vdots \\ \vdots & \ddots & \mathbf{I} & \ddots & 0 \\ \vdots & & \ddots & \ddots & -A_K^* \\ 0 & \dots & \dots & 0 & \mathbf{I} \end{pmatrix}.$$

Solving an upper triangular system like $\partial_1 c(\mathbf{s}(\mathbf{w}), \mathbf{w})^* \mathbf{r} = \mathbf{u}$ can then be done efficiently by **backsubstitution**. Starting from the last adjoint state $\mathbf{r}_K = \mathbf{u}$, we can compute each adjoint state \mathbf{r}_k from that computed at $k + 1$. Namely, for $k \in (K - 1, \dots, 1)$, we have

$$\mathbf{r}_k - A_{k+1}^* \mathbf{r}_{k+1} = \mathbf{0} \iff \mathbf{r}_k = \partial f_{k+1}(\mathbf{s}_k, \mathbf{w}_{k+1})^* \mathbf{r}_{k+1}.$$

The VJPs with respect to the parameters are then obtained by

$$\partial_2 c(\mathbf{s}(\mathbf{w}), \mathbf{w})^* \mathbf{r} = \begin{pmatrix} \partial_2 f_1(\mathbf{x}, \mathbf{w}_1)^* \mathbf{r}_1 \\ \partial_2 f_2(\mathbf{s}_1(\mathbf{w}), \mathbf{w}_2)^* \mathbf{r}_2 \\ \vdots \\ \partial_2 f_K(\mathbf{s}_1(\mathbf{w}), \mathbf{w}_K)^* \mathbf{r}_K \end{pmatrix},$$

recovering reverse-mode autodiff.

The Lagrangian perspective of backpropagation for networks with separate parameters $\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_K)$ is well-known; see for instance LeCun (1988) or Recht (2016). The Lagrangian perspective of backpropagation through time (Werbos, 1990) for networks with shared parameter \mathbf{w} is discussed for instance by Franceschi *et al.* (2017). Our exposition uses the adjoint state method, which can itself be proved either using the method of Lagrange multipliers (Section 11.4.3) or by the implicit function theorem (Section 11.4.4), combined with backsubstitution for solving the upper-triangular linear system. Past works often minimize over \mathbf{w} but we do not require this, as gradients are not necessarily used for optimization. Our exposition also supports computing the VJP of any vector-valued function f , while existing works derive the gradient of a scalar-valued loss function.

11.5 Inverse function theorem

11.5.1 Differentiating inverse functions

In some cases (see for instance Section 12.4.4), it is useful to compute the Jacobian of an inverse function f^{-1} . The inverse function theorem below allows us to relate the Jacobian of f^{-1} with the Jacobian of f .

Theorem 11.5 (Inverse function theorem). Assume $f: \mathcal{W} \rightarrow \mathcal{W}$ is continuously differentiable with invertible Jacobian $\partial f(\mathbf{w}_0)$ at \mathbf{w}_0 . Then f is bijective from a neighborhood of \mathbf{w}_0 to a neighborhood of $f(\mathbf{w}_0)$. Moreover, the inverse f^{-1} is continuously differentiable near $\boldsymbol{\omega}_0 = f(\mathbf{w}_0)$ and the Jacobian of the inverse $\partial f^{-1}(\boldsymbol{\omega})$ is

$$\partial f(\mathbf{w}) \partial f^{-1}(\boldsymbol{\omega}) = I \Leftrightarrow \partial f^{-1}(\boldsymbol{\omega}) = (\partial f(\mathbf{w}))^{-1},$$

| with $\mathbf{w} = f^{-1}(\boldsymbol{\omega})$.

11.5.2 Link with the implicit function theorem

The inverse function theorem can be used to prove the implicit function theorem; see proof of Theorem 11.4. Conversely, recall that, in order to use the implicit function theorem, we need to choose a root objective $F: \mathcal{W} \times \Lambda \rightarrow \mathcal{W}$. If we set $\mathcal{W} = \Lambda = \mathbb{R}^Q$ and $F(\mathbf{w}, \boldsymbol{\omega}) = f(\mathbf{w}) - \boldsymbol{\omega}$, with $f: \mathbb{R}^Q \rightarrow \mathbb{R}^Q$, then we have that the root $\mathbf{w}^*(\boldsymbol{\omega})$ satisfying $F(\mathbf{w}^*(\boldsymbol{\omega}), \boldsymbol{\omega}) = \mathbf{0}$ is exactly $\mathbf{w}^*(\boldsymbol{\omega}) = f^{-1}(\boldsymbol{\omega})$. Moreover, $\partial_1 F(\mathbf{w}, \boldsymbol{\omega}) = \partial f(\mathbf{w})$ and $\partial_2 F(\mathbf{w}, \boldsymbol{\omega}) = -I$. By applying the implicit function theorem with this F , we indeed recover the inverse function theorem.

11.5.3 Proof of inverse function theorem

We first give a proof of the formula assuming that f^{-1} is well-defined and continuously differentiable in a neighborhood of $f(\mathbf{w}_0)$. In that case, we have for any $\boldsymbol{\omega}$ in a neighborhood of $f(\mathbf{w}_0)$,

$$f \circ f^{-1}(\boldsymbol{\omega}) = \boldsymbol{\omega}.$$

Differentiating both sides w.r.t. $\boldsymbol{\omega}$, we get

$$\partial f(f^{-1}(\boldsymbol{\omega})) \partial f^{-1}(\boldsymbol{\omega}) = I,$$

where I is the identity function in \mathbb{R}^Q . In particular, for $\mathbf{w} = f^{-1}(\boldsymbol{\omega})$ we recover the formula presented in the statement.

Now, it remains to show that invertibility of the JVP ensures that the function is invertible in a neighborhood of $f(\mathbf{w}_0)$ and that the inverse is continuously differentiable. For that, denote $\mathbf{l} = \partial f(\mathbf{w}_0)$ such that \mathbf{l}^{-1} is well-defined by definition. f is invertible with continuously differentiable inverse, if and only if $\mathbf{l}^{-1}(f(\mathbf{w})) - f(\mathbf{w}_0)$ is invertible with continuously differentiable inverse. So without loss of generality, we consider $\partial f(\mathbf{w}_0) = I$, $f(\mathbf{w}_0) = 0$, $\mathbf{w}_0 = 0$.

As f is continuously differentiable, there exists a neighborhood $N = \{\mathbf{w} : \|\mathbf{w} - \mathbf{w}_0\|_2 \leq \delta\}$ on which we have $\|\partial f(\mathbf{w}) - I\|_2 \leq 1/2$. In this neighborhood, the function $g(\mathbf{w}) = f(\mathbf{w}) - \mathbf{w}$ is contractive by the mean value theorem with contraction factor $1/2$. For any $\boldsymbol{\omega}$ such that

$\|\omega - f(\mathbf{w}_0)\|_2 \leq \delta/2$, the sequence $\mathbf{w}_{k+1} = \mathbf{w}_k - f(\mathbf{w}_k) - \omega'$ remains in \mathcal{N} and converges (since it is a Cauchy sequence by the contraction of g) to a unique fixed point \mathbf{w}_∞ satisfying $\mathbf{w}_\infty = \mathbf{w}_\infty - f(\mathbf{w}_\infty) - \omega \iff f(\mathbf{w}_\infty) = \omega$. This shows the existence of the inverse in the neighborhood $M = \{\omega : \|\omega - \omega_0\|_2 \leq \delta/2\}$ of $\omega_0 = f(\mathbf{w}_0)$ onto N .

We tackle now the differentiability (hence the continuity) of f^{-1} . For any ω in the neighborhood of ω_0 with inverse $\mathbf{w} := f^{-1}(\omega) \in N$, the JVP of f at \mathbf{w} satisfies by assumption $\|\partial f(\mathbf{w}) - \mathbf{I}\|_2 \leq 1/2$. Hence, $\mathbf{a} = \partial f(\mathbf{w}) - \mathbf{I}$ defines a convergent series $\mathbf{b} = \sum_{k=0}^{+\infty} \mathbf{a}^k$ and one verifies easily that $\mathbf{b} = \partial f(\mathbf{w})^{-1}$, that is $\partial f(\mathbf{w})$ is invertible and $\|\partial f(\mathbf{w})^{-1}\|_2 \leq 2$. To compute the JVP of the inverse, we consider then $\partial f(\mathbf{w})^{-1}$ as the candidate JVP and examine

$$\frac{\|f^{-1}(\omega + \eta) - f^{-1}(\omega) - (\partial f(\mathbf{w}))^{-1}\eta\|_2}{\|\eta\|_2}.$$

Denote then \mathbf{v} such that $f^{-1}(\omega + \eta) = \mathbf{w} + \mathbf{v}$. As $g(\mathbf{w}) = f(\mathbf{w}) - \mathbf{w}$ is $1/2$ -contractive in N , we have $\|\mathbf{v} - \eta\|_2 = \|g(\mathbf{w} + \mathbf{v}) - g(\mathbf{w})\|_2 \leq 1/2\|\mathbf{v}\|_2$. So $\|\mathbf{v}\|_2 \geq \|\eta\|_2/2$. We then get

$$\begin{aligned} & \frac{\|f^{-1}(\omega + \eta) - f^{-1}(\omega) - (\partial f(\mathbf{w}))^{-1}\eta\|_2}{\|\eta\|_2} \\ &= \frac{\|\mathbf{v} - (\partial f(\mathbf{w}))^{-1}(f(\mathbf{w} + \mathbf{v}) - f(\mathbf{w}))\|_2}{\|\eta\|_2} \\ &\leq 4 \frac{\|f(\mathbf{w} + \mathbf{v}) - f(\mathbf{w}) - \partial f(\mathbf{w})\mathbf{v}\|_2}{\|\mathbf{v}\|_2} \end{aligned}$$

As $\|\eta\|_2 \rightarrow 0$, we have $\|\mathbf{v}\|_2 \rightarrow 0$ and so $\|f(\mathbf{w} + \mathbf{v}) - f(\mathbf{w}) - \partial f(\mathbf{w})\mathbf{v}\|_2 / \|\mathbf{v}\|_2 \rightarrow 0$. Hence, f^{-1} is differentiable with JVP $\partial f^{-1}(\omega) = (\partial f(\mathbf{w}))^{-1} = (\partial f(f^{-1}(\omega)))^{-1}$. This shows that f^{-1} is continuous and so $\partial f^{-1}(\omega)$ is continuous as a composition of continuous functions.

11.6 Summary

- Implicit functions are functions that cannot be decomposed into elementary operations and for which autodiff can therefore not be directly applied. Examples are optimization problems and nonlinear equations.

- Envelope theorems can be used for differentiating through the min or max value (not solution) of a function.
- More generally, the implicit function theorem allows us to differentiate through implicit functions. It gives conditions for the existence of derivatives and how to obtain them.
- The adjoint state method can be used to obtain the gradient of the composition of an explicit function and of an implicit function, specified by equality constraints. It can be used to prove the correctness of reverse-mode autodiff.
- The inverse function theorem can be used to differentiate function inverses.
- In a sense, the implicit function theorem can be thought as the mother theorem, as it can be used to prove envelope theorems, the adjoint state method and the inverse function theorem.

12

Differentiating through integration

In this chapter, we study how to differentiate through integrals, with a focus on expectations and solutions of ordinary differential equations.

12.1 Differentiation under the integral sign

Given two Euclidean spaces Θ and \mathcal{Y} , and a function $f : \Theta \times \mathcal{Y} \rightarrow \mathbb{R}$, we often want to differentiate an integral of the form

$$F(\boldsymbol{\theta}) := \int_{\mathcal{Y}} f(\boldsymbol{\theta}, \mathbf{y}) d\mathbf{y}.$$

Provided that we can swap integration and differentiation, we have

$$\nabla F(\boldsymbol{\theta}) = \int_{\mathcal{Y}} \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}, \mathbf{y}) d\mathbf{y}.$$

The conditions enabling us to do so are best examined in the context of measure theory. We refer the reader to e.g. (Cohn, 2013) for a course on measure theory and Flanders (1973) for an in-depth study of the differentiation under the integral sign. Briefly, if $\Theta = \mathcal{Y} = \mathbb{R}$, the following conditions are sufficient.

1. f is measurable in both its arguments, and $f(\boldsymbol{\theta}, \cdot)$ is integrable for almost all $\boldsymbol{\theta} \in \Theta$ fixed,

2. $f(\cdot, y)$ is absolutely continuous for almost all $y \in \mathcal{Y}$, that is, there exists an integrable function $g(\cdot, y)$ such that $f(\theta, y) = f(\theta_0, y) + \int_{\theta_0}^{\theta} g(\tau, y) d\tau$,
3. $\partial_1 f(\theta, y)$ (which exists almost everywhere if $f(\cdot, y)$ is absolutely continuous), is locally integrable, that is, for any closed interval $[\theta_0, \theta_1]$, the integral $\int_{\theta_0}^{\theta_1} \int |\partial_1 f(\theta, y)| dy d\theta$ is finite.

Any differentiable function $f : \Theta \times \mathcal{Y} \rightarrow \mathbb{R}$ is absolutely continuous. However, the conditions also hold if f is just absolutely continuous, that is, if $f(\cdot, y)$ is differentiable for almost all y . This weaker assumption can be used to smooth out differentiable almost-everywhere functions, such as the ReLu, as we study in Section 14.4.

12.2 Differentiating through expectations

A special case of differentiating through integrals is differentiating through expectations. We can distinguish between two cases, depending on whether the parameters θ we wish to differentiate are involved in the distribution or in the function, whose expectation we compute.

12.2.1 Parameter-independent distributions

We first consider expectations of the form

$$F(\theta) := \mathbb{E}_{Y \sim p}[g(Y, \theta)] = \int_{\mathcal{Y}} g(\mathbf{y}, \theta) p(\mathbf{y}) d\mathbf{y},$$

for a random variable $Y \in \mathcal{Y} \subseteq \mathbb{R}^M$, distributed according to a distribution p , and a function $g : \mathcal{Y} \times \Theta \rightarrow \mathbb{R}$. Importantly, the distribution is independent of the parameters θ . Under mild conditions recalled in Section 12.1, we can swap differentiation and integration to obtain

$$\begin{aligned} \nabla F(\theta) &= \nabla_{\theta} \int_{\mathcal{Y}} g(\mathbf{y}, \theta) p(\mathbf{y}) d\mathbf{y} \\ &= \int_{\mathcal{Y}} \nabla_{\theta} g(\mathbf{y}, \theta) p(\mathbf{y}) d\mathbf{y} \\ &= \mathbb{E}_{Y \sim p}[\nabla_{\theta} g(Y, \theta)]. \end{aligned}$$

Generally, the expectation cannot be computed in closed form. However, provided that we can sample from p , we can define a Monte-Carlo estimator of the value

$$\hat{F}_N(\boldsymbol{\theta}) := \frac{1}{N} \sum_{i=1}^N g(Y_i, \boldsymbol{\theta})$$

and of the gradient

$$\nabla \hat{F}_N(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \nabla_{\boldsymbol{\theta}} g(Y_i, \boldsymbol{\theta}),$$

for N i.i.d. samples Y_1, \dots, Y_N from p . These estimators are unbiased, meaning that $\mathbb{E}[\hat{F}_N(\boldsymbol{\theta})] = F(\boldsymbol{\theta})$ and $\mathbb{E}[\nabla \hat{F}_N(\boldsymbol{\theta})] = \nabla F(\boldsymbol{\theta})$, and converge to the true quantity as $N \rightarrow +\infty$. This suggests a simple implementation in an autodiff framework of the approximation of $\nabla F(\boldsymbol{\theta})$:

1. Sample y_1, \dots, y_n from p .
2. Compute $\hat{F}_N(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n g(y_i, \boldsymbol{\theta})$.
3. Compute the gradient $\nabla \hat{F}_N(\boldsymbol{\theta})$ by automatic differentiation.

Computing higher order derivatives follow the same principle: to get an approximation of $\nabla^2 F(\boldsymbol{\theta})$, we can simply compute $\nabla^2 \hat{F}_N(\boldsymbol{\theta})$ by autodiff. As such, the implementation delineated above is akin to the “discretize-then-optimize” approach used to differentiate through the solution of an ODE (Section 12.6): we implement an approximation of the objective and simply call autodiff on it.

12.2.2 Parameter-dependent distributions

A more challenging case arises when the distribution depends on the parameters $\boldsymbol{\theta}$:

$$E(\boldsymbol{\theta}) := \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}}[g(Y)] = \int_{\mathcal{Y}} g(\mathbf{y}) p_{\boldsymbol{\theta}}(\mathbf{y}) d\mathbf{y},$$

where $Y \in \mathcal{Y} \subseteq \mathbb{R}^M$ is a random variable, distributed according to a distribution $p_{\boldsymbol{\theta}}$ parameterized by $\boldsymbol{\theta} \in \Theta$ and where $g: \mathcal{Y} \rightarrow \mathbb{R}$ is, depending on the setting, potentially a blackbox function (i.e., we do not

have access to its gradients). Typically, $\boldsymbol{\theta} \in \Theta$ could be parameters we wish to estimate, or it could indirectly be generated by $\boldsymbol{\theta} = f(\mathbf{x}, \mathbf{w}) \in \Theta$, where f is a neural network with parameters $\mathbf{w} \in \mathcal{W}$ we wish to estimate. The main difficulty in computing $\nabla E(\boldsymbol{\theta})$ stems from the fact that $\boldsymbol{\theta}$ are the parameters of the distribution $p_{\boldsymbol{\theta}}$. Estimating an expectation $E(\boldsymbol{\theta}) = \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}}[g(Y)]$ using Monte-Carlo estimation requires us to sample from $p_{\boldsymbol{\theta}}$. However, it is not clear how to differentiate E w.r.t. $\boldsymbol{\theta}$ if $\boldsymbol{\theta}$ is involved in the sampling process.

Continuous case

When \mathcal{Y} is a continuous set (that is, $p_{\boldsymbol{\theta}}(\mathbf{y})$ is a probability density function), we can rewrite $E(\boldsymbol{\theta})$ as

$$E(\boldsymbol{\theta}) = \int_{\mathcal{Y}} p_{\boldsymbol{\theta}}(\mathbf{y}) g(\mathbf{y}) d\mathbf{y}.$$

Provided that we can swap integration and differentiation (see Section 12.1), we then have

$$\begin{aligned} \nabla E(\boldsymbol{\theta}) &= \nabla_{\boldsymbol{\theta}} \int_{\mathcal{Y}} p_{\boldsymbol{\theta}}(\mathbf{y}) g(\mathbf{y}) d\mathbf{y} \\ &= \int_{\mathcal{Y}} \nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\mathbf{y}) g(\mathbf{y}) d\mathbf{y}. \end{aligned}$$

Unfortunately, this integral is not an expectation and it could be intractable in general.

Discrete case

When \mathcal{Y} is a discrete set (that is, $p_{\boldsymbol{\theta}}(\mathbf{y})$ is a probability mass function), we can rewrite $E(\boldsymbol{\theta})$ as

$$E(\boldsymbol{\theta}) = \sum_{\mathbf{y} \in \mathcal{Y}} p_{\boldsymbol{\theta}}(\mathbf{y}) g(\mathbf{y}).$$

We then obtain

$$\nabla E(\boldsymbol{\theta}) = \sum_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{y}) \nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\mathbf{y}).$$

Again $\nabla E(\boldsymbol{\theta})$ is not an expectation. We therefore cannot use Monte-Carlo estimation to estimate the gradient. Instead, we can compute it

by brute force, i.e., by summing over all possible $\mathbf{y} \in \mathcal{Y}$. However, this is clearly only computationally tractable if $|\mathcal{Y}|$ is small or if p_{θ} is designed to have sparse support, i.e., so that the set $\{\mathbf{y} \in \mathcal{Y}: p_{\theta}(\mathbf{y}) \neq 0\}$ is small. Moreover, even if these conditions hold, summing over \mathbf{y} could be problematic if $g(\mathbf{y})$ is expensive to compute. Therefore, exact gradients are seldom used in practice.

In Sections 12.3 and 12.4, we review the score function and pathwise gradient estimators, to (approximately) compute $\nabla E(\theta)$, allowing us to optimize θ (or \mathbf{w} using the chain rule) by gradient-based algorithms.

12.2.3 Application to expected loss functions

Differentiating through expectations is particularly useful when working with expected loss functions of the form

$$L(\theta; \mathbf{y}) := \mathbb{E}_{\hat{\mathbf{Y}} \sim p_{\theta}}[\ell(\hat{\mathbf{Y}}, \mathbf{y})],$$

where \mathbf{y} is some ground truth. Equivalently, we can set $\ell = -r$, where r is a **reward function**. As we shall see, the score function estimator will support a discrete loss function $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, while the pathwise gradient estimator will require a differentiable loss function $\ell: \mathbb{R}^M \times \mathcal{Y} \rightarrow \mathbb{R}$. Intuitively, $L(\theta; \mathbf{y})$ will be low if p_{θ} assigns high probability to predictions $\hat{\mathbf{y}}$ with low loss value $\ell(\hat{\mathbf{y}}, \mathbf{y})$.

In the classification setting, where $\mathcal{Y} = [M]$, p_{θ} is often chosen to be the **Gibbs distribution**, which is a categorical distribution induced by a softmax

$$p_{\theta}(y) := \frac{\exp(\theta_y)}{\sum_{i \in [M]} \exp(\theta_i)} = [\text{softmax}(\theta)]_y \in (0, 1),$$

where $\theta_y := f(\mathbf{x}, y, \mathbf{w}) \in \mathbb{R}$ are logits produced by a neural network f . More generally, in the structured prediction setting, where $\mathcal{Y} \subseteq \mathbb{R}^M$ but $|\mathcal{Y}| \gg M$, we often use the distribution

$$p_{\theta}(\mathbf{y}) := \frac{\exp(\langle \phi(\mathbf{y}), \theta \rangle)}{\sum_{\mathbf{y}' \in \mathcal{Y}} \exp(\langle \phi(\mathbf{y}'), \theta \rangle)},$$

where $\theta = f(\mathbf{x}, \mathbf{w}) \in \mathbb{R}^M$.

Given a distribution ρ over $\mathcal{X} \times \mathcal{Y}$, we then want to minimize the expected loss function, also known as **risk**,

$$R(\mathbf{w}) := \mathbb{E}_{(X,Y) \sim \rho}[L(f(X, \mathbf{w}); Y)].$$

Typically, minimizing $R(\mathbf{w})$ is done through some form of gradient descent, which requires us to be able to compute

$$\begin{aligned} \nabla R(\mathbf{w}) &= \mathbb{E}_{(X,Y) \sim \rho}[\nabla_{\mathbf{w}} L(f(X, \mathbf{w}); Y)] \\ &= \mathbb{E}_{(X,Y) \sim \rho}[\partial_2 f(\mathbf{x}, \mathbf{w})^* \nabla L(f(X, \mathbf{w}); Y)]. \end{aligned}$$

Computing $\nabla R(\mathbf{w})$ therefore boils down to computing the gradient of $L(\boldsymbol{\theta}; \mathbf{y})$, which is the gradient of an expectation.

12.2.4 Application to experimental design

In experimental design, we wish to minimize a function $g(\boldsymbol{\lambda})$, which we assume costly to evaluate. As an example, evaluating $g(\boldsymbol{\lambda})$ could require us to run a scientific experiment with parameters $\boldsymbol{\lambda} \in \mathbb{R}^Q$. As another example, in hyperparameter optimization, evaluating $g(\boldsymbol{\lambda})$ would require us to run a learning algorithm with hyperparameters $\boldsymbol{\lambda} \in \mathbb{R}^Q$. Instead of solving the problem $\arg \min_{\boldsymbol{\lambda} \in \mathbb{R}^Q} g(\boldsymbol{\lambda})$, we can lift the problem to probability distributions and solve $\arg \min_{\boldsymbol{\theta} \in \mathbb{R}^M} E(\boldsymbol{\theta})$, where $E(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\lambda} \sim p_{\boldsymbol{\theta}}}[g(\boldsymbol{\lambda})]$. This requires the probability distribution $p_{\boldsymbol{\theta}}$ to assign high probability to $\boldsymbol{\lambda}$ values that achieve small $g(\boldsymbol{\lambda})$ value. Solving this problem by stochastic gradient descent requires us to be able to compute estimates of $\nabla E(\boldsymbol{\theta})$. This can be done for instance with SFE explained in Section 12.3, which does not require gradients of g , unlike implicit differentiation explained in Chapter 11. This approach also requires us to choose a distribution $p_{\boldsymbol{\theta}}$ over $\boldsymbol{\lambda}$. For continuous hyperparameters, a natural choice would be the normal distribution $\boldsymbol{\lambda} \sim \text{Normal}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, setting $\boldsymbol{\theta} = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Once we obtained $\boldsymbol{\theta}$ by minimizing $E(\boldsymbol{\theta})$, we need a way to recover $\boldsymbol{\lambda}$. This can be done for example by choosing the mode of the distribution, i.e., $\arg \max_{\boldsymbol{\lambda} \in \mathbb{R}^Q} p_{\boldsymbol{\theta}}(\boldsymbol{\lambda})$, or the mean of the distribution $\mathbb{E}_{\boldsymbol{\lambda} \sim p_{\boldsymbol{\theta}}(\boldsymbol{\lambda})}[\boldsymbol{\lambda}]$. Of course, in the case of the normal distribution, they coincide.

12.3 Score function estimators, REINFORCE

12.3.1 Scalar-valued functions

The key idea of the **score function estimator** (SFE), also known as REINFORCE, is to rewrite $\nabla E(\boldsymbol{\theta})$ as an expectation. The estimator is based on the **logarithmic derivative identity**

$$\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{y}) = \frac{\nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\mathbf{y})}{p_{\boldsymbol{\theta}}(\mathbf{y})} \iff \nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\mathbf{y}) = p_{\boldsymbol{\theta}}(\mathbf{y}) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{y}).$$

Using this identity, we obtain the following gradient estimator.

Proposition 12.1 (SFE for scalar-valued functions). Given a family of distributions $p_{\boldsymbol{\theta}}$ on \mathcal{Y} , for $\boldsymbol{\theta} \in \Theta$, define

$$E(\boldsymbol{\theta}) := \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}} [g(Y)] = \int_{\mathcal{Y}} p_{\boldsymbol{\theta}}(\mathbf{y}) g(\mathbf{y}) d\mathbf{y},$$

where $Y \in \mathcal{Y} \subseteq \mathbb{R}^M$ and $g: \mathcal{Y} \rightarrow \mathbb{R}$. Then,

$$\nabla E(\boldsymbol{\theta}) = \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}} [g(Y) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(Y)].$$

Proof.

$$\begin{aligned} \nabla E(\boldsymbol{\theta}) &= \int_{\mathcal{Y}} \nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\mathbf{y}) g(\mathbf{y}) d\mathbf{y} \\ &= \int_{\mathcal{Y}} p_{\boldsymbol{\theta}}(\mathbf{y}) g(\mathbf{y}) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{y}) d\mathbf{y} \\ &= \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}} [g(Y) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(Y)]. \end{aligned}$$

□

The gradient of the log-PDF w.r.t. $\boldsymbol{\theta}$, $\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{y})$, is known as the **score function**, hence the estimator name. SFE is suitable when two requirements are met: it is easy to sample from $p_{\boldsymbol{\theta}}$ and the score function is available in closed form. Since the SFE gradient is an expectation, we can use Monte-Carlo estimation to compute an unbiased estimator of $\nabla E(\boldsymbol{\theta})$:

$$\nabla E(\boldsymbol{\theta}) \approx \hat{\gamma}_N(\boldsymbol{\theta}) := \frac{1}{N} \sum_{i=1}^N g(Y_i) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(Y_i), \quad (12.1)$$

where Y_1, \dots, Y_N are sampled from p_θ .

Interestingly, the gradient of g is not needed in this estimator. Therefore, there is no differentiability assumption about g . This is why SFE is useful when g is a discrete loss function or more generally a blackbox function.

Example 12.1 (SFE with a language model). In a language model, the probability of a sentence $\mathbf{y} = (y_1, \dots, y_L)$ is typically factored using the chain rule of probability (see Section 10.1)

$$p_\theta(\mathbf{y}) := p_\theta(y_1)p_\theta(y_2|y_1) \dots p_\theta(y_L|y_1, \dots, y_{L-1}),$$

where p_θ is modeled using a transformer or RNN. Note that the probabilities are normalized by construction, so there is no need for an explicit normalization constant. Thanks to this factorization, it is easy to sample from p_θ using ancestral sampling (see Section 10.5.3) and the log-probability enjoys the simple expression

$$\begin{aligned} \nabla_\theta \log p_\theta(\mathbf{y}) &= \nabla_\theta \log p_\theta(y_1) + \nabla_\theta \log p_\theta(y_2|y_1) + \dots \\ &\quad + \nabla_\theta \log p_\theta(y_L|y_1, \dots, y_{L-1}). \end{aligned}$$

This gradient is easy to compute, since the token-wise distributions $p_\theta(y_j|y_1, \dots, y_{j-1})$ are typically defined using a softargmax. We can therefore easily compute $\nabla E(\theta)$ under p_θ using SFE. This is for instance useful to optimize an expected reward, in order to finetune or align a language model (Ziegler *et al.*, 2019).

Another example when $\nabla_\theta p_\theta(\mathbf{y})$ is available in closed form is in the context of reinforcement learning, where $p_\theta(\mathbf{y})$ is a Markov Decision Process (MDP) and is called the policy. Applying the SFE leads to the (vanilla) policy gradient method (Sutton *et al.*, 1999) and can then be used to compute the gradient of an expected cumulative reward. However, SFE is more problematic when used with the Gibbs distribution, due to the explicit normalization constant.

Example 12.2 (SFE with a Gibbs distribution). The Gibbs distribu-

tion is parameterized, for $\boldsymbol{\theta} \in \mathbb{R}^{\mathcal{Y}}$,

$$p_{\boldsymbol{\theta}}(y) := \exp(\theta_y/\gamma - A(\boldsymbol{\theta})) = \exp(\theta_y/\gamma) / \exp(A(\boldsymbol{\theta}))$$

where we defined the log-partition function

$$A(\boldsymbol{\theta}) := \log \sum_{y \in \mathcal{Y}} \exp(\theta_y/\gamma).$$

A typical parametrization is $\theta_y = f(\mathbf{x}, y, \mathbf{w})$ with f the output of network on a sample \mathbf{x} with parameters \mathbf{w} . We then have

$$\log p_{\boldsymbol{\theta}}(y) = \theta_y/\gamma - A(\boldsymbol{\theta}),$$

so that

$$\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(y) = \mathbf{e}_y/\gamma - \nabla A(\boldsymbol{\theta}).$$

We therefore see that $\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(y)$ crucially depends on $\nabla A(\boldsymbol{\theta})$, the gradient of the log-partition. This gradient is available for some structured sets \mathcal{Y} , see e.g. (Mensch and Blondel, 2018), but not in general.

As another example, we apply SFE in Section 14.4 to derive the gradient of perturbed functions.

Differentiating through both the distribution and the function

Suppose both the distribution and the function now depend on $\boldsymbol{\theta}$. When g is scalar-valued and differentiable w.r.t. $\boldsymbol{\theta}$, we want to differentiate

$$E(\boldsymbol{\theta}) := \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}} [g(Y, \boldsymbol{\theta})].$$

Using the product rule, we obtain

$$\nabla E(\boldsymbol{\theta}) = \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}} [g(Y, \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(Y)] + \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}} [\nabla_{\boldsymbol{\theta}} g(Y, \boldsymbol{\theta})].$$

Differentiating through joint distributions

Suppose we now want to differentiate through

$$E(\boldsymbol{\theta}) := \mathbb{E}_{Y_1 \sim p_{\boldsymbol{\theta}}, Y_2 \sim q_{\boldsymbol{\theta}}} [g(Y_1, Y_2)].$$

The gradient is then given by

$$\nabla E(\boldsymbol{\theta}) = \mathbb{E}_{Y_1 \sim p_{\boldsymbol{\theta}}, Y_2 \sim q_{\boldsymbol{\theta}}}[(\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(Y_1) + \nabla \log q_{\boldsymbol{\theta}}(Y_2))g(Y_1, Y_2)],$$

which is easily seen by applying Proposition 12.1 on the joint distribution $\rho_{\boldsymbol{\theta}} := p_{\boldsymbol{\theta}} \cdot q_{\boldsymbol{\theta}}$. The extension to more than two variables is straightforward.

12.3.2 Variance reduction

Bias and variance

Recall the definition of $\hat{\gamma}_N$ in Eq. (12.1). SFE is an **unbiased** estimator, meaning that

$$\nabla E(\boldsymbol{\theta}) = \mathbb{E}[\hat{\gamma}_N(\boldsymbol{\theta})],$$

where the expectation is taken with respect to the N samples drawn. Since the gradient is vector-valued, we need to define a scalar-valued notion of variance. We do so by using the squared Euclidean distance in the usual variance definition to define

$$\begin{aligned} \mathbb{V}[\hat{\gamma}_N(\boldsymbol{\theta})] &:= \mathbb{E}[\|\hat{\gamma}_N(\boldsymbol{\theta}) - \nabla E(\boldsymbol{\theta})\|_2^2] \\ &= \mathbb{E}[\|\hat{\gamma}_N(\boldsymbol{\theta})\|_2^2] - \|\nabla E(\boldsymbol{\theta})\|_2^2. \end{aligned}$$

The variance naturally goes to zero as $N \rightarrow \infty$.

Baseline

SFE is known to suffer from **high variance** (Mohamed *et al.*, 2020). This means that this estimator may require us to draw many samples from the distribution $p_{\boldsymbol{\theta}}$ to work well in practice. One of the simplest variance reduction technique consists in shifting the function g with a constant β , called a **baseline**, to obtain

$$\nabla E(\boldsymbol{\theta}) = \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}}[(g(Y) - \beta)\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(Y)].$$

The reason this is still a valid estimator of $\nabla E(\boldsymbol{\theta})$ stems from

$$\begin{aligned} \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}}[\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(Y)] &= \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}} \left[\frac{\nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(Y)}{p_{\boldsymbol{\theta}}(Y)} \right] \\ &= \nabla_{\boldsymbol{\theta}} \mathbb{E}_{Y \sim p_{\boldsymbol{\theta}}}[\mathbf{1}] \\ &= \nabla_{\boldsymbol{\theta}} \mathbf{1} \\ &= \mathbf{0}, \end{aligned}$$

for any valid distribution p_{θ} . The baseline β is often set to the running average of past values of the function g , though it is neither optimal nor does it guarantee to lower the variance (Mohamed *et al.*, 2020).

Control variates

Another general technique are **control variates**. Let us denote the expectation of a function $h: \mathbb{R}^M \rightarrow \mathbb{R}$ under the distribution p_{θ} as

$$H(\theta) := \mathbb{E}_{Y \sim p_{\theta}}[h(Y)].$$

Suppose that $H(\theta)$ and its gradient $\nabla H(\theta)$ are known in closed form. Then, for any $\gamma \geq 0$, we clearly have

$$\begin{aligned} E(\theta) &= \mathbb{E}_{Y \sim p_{\theta}}[g(Y)] \\ &= \mathbb{E}_{Y \sim p_{\theta}}[g(Y) - \gamma(h(Y) - H(\theta))] \\ &= \mathbb{E}_{Y \sim p_{\theta}}[g(Y) - \gamma h(Y)] + \gamma H(\theta) \end{aligned}$$

and therefore

$$\nabla E(\theta) = \nabla_{\theta} \mathbb{E}_{Y \sim p_{\theta}}[g(Y) - \gamma h(Y)] + \gamma \nabla H(\theta).$$

Applying SFE, we then obtain

$$\nabla E(\theta) = \mathbb{E}_{Y \sim p_{\theta}}[(g(Y) - \gamma h(Y)) \nabla_{\theta} \log p_{\theta}(Y)] + \gamma \nabla H(\theta).$$

Examples of h include a bound on f or a second-order Taylor expansion of f , assuming that these approximations are easier to integrate than f (Mohamed *et al.*, 2020).

12.3.3 Vector-valued functions

It is straightforward to extend the SFE to vector-valued functions.

Proposition 12.2 (SFE for vector-valued functions). Given a family of distributions p_{θ} on \mathcal{Y} , for $\theta \in \Theta$, define

$$E(\theta) := \mathbb{E}_{Y \sim p_{\theta}}[g(Y)] = \int_{\mathcal{Y}} p_{\theta}(\mathbf{y}) g(\mathbf{y}) d\mathbf{y},$$

where $Y \in \mathcal{Y}$, $g: \mathcal{Y} \rightarrow \mathcal{G}$. The JVP of E at $\theta \in \Theta$ along $\mathbf{v} \in \Theta$ is

$$\partial E(\theta) \mathbf{v} = \mathbb{E}_{Y \sim p_{\theta}}[\langle \nabla_{\theta} \log p_{\theta}(Y), \mathbf{v} \rangle g(Y)] \in \mathcal{G}$$

and the VJP of E at $\theta \in \Theta$ along $u \in \mathcal{G}$ is

$$\partial E(\theta)^* u = \mathbb{E}_{Y \sim p_\theta} [\nabla_\theta \log p_\theta(Y) \langle u, g(Y) \rangle] \in \Theta$$

The Jacobian of E at $\theta \in \Theta$ can then be written as

$$\partial E(\theta) = \mathbb{E}_{Y \sim p_\theta} [g(Y) \otimes \nabla_\theta \log p_\theta(Y)],$$

where \otimes denote the outer product.

Proof. The VJP of E at $\theta \in \Theta$ along $u \in \Theta$ amounts to compute the gradient of the scalar function

$$\langle E(\theta), u \rangle = \mathbb{E}_{Y \sim p_\theta} [\langle g(Y), u \rangle]$$

The expression of the VJP follows by using the SFE on the scalar valued integrand $\langle g(Y), u \rangle$. The JVP is obtained as the adjoint operator of the VJP and the Jacobian follows. \square

Differentiating through both the distribution and the function

If θ now influences both the distribution and the function,

$$E(\theta) := \mathbb{E}_{Y \sim p_\theta} [g(Y, \theta)],$$

then, we obtain

$$\partial E(\theta) = \mathbb{E}_{Y \sim p_\theta} [g(Y, \theta) \otimes \nabla_\theta \log p_\theta(Y)] + \mathbb{E}_{Y \sim p_\theta} [\partial_\theta g(Y, \theta)].$$

12.3.4 Second derivatives

Using the previous subsection with $g(y, \theta) = g(y) \nabla_\theta \log p_\theta(\theta)$, we easily obtain an estimator of the Hessian.

Proposition 12.3 (SFE for the Hessian). Let us define the scalar-valued function $E(\theta) := \mathbb{E}_{Y \sim p_\theta} [g(Y)]$. Then,

$$\begin{aligned} \nabla^2 E(\theta) = & \mathbb{E}_{Y \sim p_\theta} [g(Y) \nabla_\theta \log p_\theta(Y) \otimes \nabla_\theta \log p_\theta(Y)] + \\ & \mathbb{E}_{Y \sim p_\theta} [g(Y) \nabla_\theta^2 \log p_\theta(Y)]. \end{aligned}$$

This can also be derived using the second-order log-derivative

$$\nabla_{\theta}^2 \log p_{\theta}(\mathbf{y}) = \frac{1}{p_{\theta}(\mathbf{y})} \nabla_{\theta}^2 p_{\theta}(\mathbf{y}) - \frac{1}{p_{\theta}(\mathbf{y})^2} \nabla_{\theta} p_{\theta}(\mathbf{y}) \otimes \nabla_{\theta} p_{\theta}(\mathbf{y})$$

so that

$$\nabla_{\theta}^2 p_{\theta}(\mathbf{y}) = p_{\theta}(\mathbf{y}) \left[\nabla_{\theta}^2 \log p_{\theta}(\mathbf{y}) + \nabla_{\theta} \log p_{\theta}(\mathbf{y}) \otimes \nabla_{\theta} \log p_{\theta}(\mathbf{y}) \right].$$

Link with the Bartlett identities

The Bartlett identities are expressions relating the moments of the score function (gradient of the log-likelihood function). Using Proposition 12.1 with $g(\mathbf{y}) = 1$ and $\int_{\mathcal{Y}} p_{\theta}(\mathbf{y}) d\mathbf{y} = 1$, we obtain

$$\mathbb{E}_{Y \sim p_{\theta}} [\nabla_{\theta} \log p_{\theta}(Y)] = \mathbf{0}, \quad (12.2)$$

which is known as Bartlett's first identity. Similarly, using Proposition 12.3, we obtain

$$\begin{aligned} & \mathbb{E}_{Y \sim p_{\theta}} [\nabla_{\theta}^2 \log p_{\theta}(Y)] + \mathbb{E}_{Y \sim p_{\theta}} [\nabla_{\theta} \log p_{\theta}(Y) \otimes \nabla_{\theta} \log p_{\theta}(Y)] \\ &= \mathbb{E}_{Y \sim p_{\theta}} [\nabla_{\theta}^2 \log p_{\theta}(Y)] + \text{cov}[\log p_{\theta}(Y)] \\ &= \mathbf{0}, \end{aligned} \quad (12.3)$$

which is known as Bartlett's second identity.

12.4 Path gradient estimators, reparametrization trick

As we saw previously, the main difficulty in computing gradients of expectations arises when the parameters θ play a role in the distribution p_{θ} being sampled. The key idea of path gradient estimators (PGE), also known as reparametrization trick, is to rewrite the expectation in such a way that the parameters are moved from the distribution to the function, using a **change of variable**.

12.4.1 Location-scale transforms

The canonical example of path gradient estimator is differentiating through the expectation

$$E(\mu, \sigma) := \mathbb{E}_{U \sim \text{Normal}(\mu, \sigma^2)} [g(U)],$$

where $g: \mathbb{R} \rightarrow \mathbb{R}$ is a differentiable function. If we let $Z \sim \text{Normal}(0, 1)$, it is easy to check that $U = \mu + \sigma Z$. We can therefore write

$$E(\mu, \sigma) = \mathbb{E}_{Z \sim \text{Normal}(0,1)}[g(\mu + \sigma Z)].$$

The key advantage is that we can now easily compute the derivatives by mere application of the chain rule, since the parameters μ and σ are moved from the distribution to the function:

$$\begin{aligned} \frac{\partial}{\partial \mu} E(\mu, \sigma) &= \mathbb{E}_{Z \sim \text{Normal}(0,1)}[g'(\mu + \sigma Z)] \\ \frac{\partial}{\partial \sigma} E(\mu, \sigma) &= \sigma \cdot \mathbb{E}_{Z \sim \text{Normal}(0,1)}[g'(\mu + \sigma Z)]. \end{aligned}$$

The change of variable

$$U := \mu + \sigma Z \tag{12.4}$$

is called a **location-scale transform**. Such a transformation exists, not only for the normal distribution, but for **location-scale family** distributions, i.e., distributions parametrized by a location parameter μ and a scale parameter $\sigma > 0$, such that U is distributed according to a distribution in the same family as Z is distributed. Besides the normal distribution, examples of location-scale family distributions include the Cauchy distribution, the uniform distribution, the logistic distribution, the Laplace distribution, and Student's t -distribution.

We can easily relate the cumulative distribution function (CDF) and the probability density function (PDF) of Z to that of U , and vice-versa.

Proposition 12.4 (CDF and PDF of location-scale family distributions).

Let $F_Z(z) := \mathbb{P}(Z \leq z)$ and $f_Z(z) := F'_Z(z)$. If $U := \mu + \sigma Z$, then

$$\begin{aligned} F_Z(z) = F_U(\mu + \sigma z) &\iff F_U(u) = F_Z\left(\frac{u - \mu}{\sigma}\right) \\ f_Z(z) = \sigma f_U(\mu + \sigma z) &\iff f_U(u) = \frac{1}{\sigma} f_Z\left(\frac{u - \mu}{\sigma}\right). \end{aligned}$$

Proof. We have

$$\begin{aligned}
 F_Z(z) &= \mathbb{P}(Z \leq z) \\
 &= \mathbb{P}\left(\frac{U - \mu}{\sigma} \leq z\right) \\
 &= \mathbb{P}(U \leq \mu + \sigma z) \\
 &= F_U(\mu + \sigma z)
 \end{aligned}$$

and we obtain $f_Z(z)$ by differentiating $F_Z(z)$. □

12.4.2 Differentiable transforms

We can generalize the idea of path gradient estimator (PGE) to any change of variable

$$U := T(Z, \boldsymbol{\theta}),$$

where $T: \mathbb{R}^M \times \mathbb{R}^Q \rightarrow \mathbb{R}^M$ is a differentiable transformation. For example, if we gather μ and σ as $\boldsymbol{\theta} := (\mu, \sigma)$, we can write the location-scale transform as

$$U = T(Z, \boldsymbol{\theta}) = \mu + \sigma Z.$$

We can derive the path gradient estimator for any such differentiable transformation T .

Proposition 12.5 (Path gradient estimator). Let us define

$$E(\boldsymbol{\theta}) := \mathbb{E}_{U \sim p_{\boldsymbol{\theta}}}[g(U)],$$

where $U \in \mathcal{U} \subseteq \mathbb{R}^M$ and $g: \mathbb{R}^M \rightarrow \mathbb{R}$ is differentiable. Suppose there is a differentiable transformation $T: \mathbb{R}^M \times \mathbb{R}^Q \rightarrow \mathbb{R}^M$ such that if $Z \sim p$ (where p does not depend on $\boldsymbol{\theta}$) and $U := T(Z, \boldsymbol{\theta})$, then $U \sim p_{\boldsymbol{\theta}}$. Then, we have

$$E(\boldsymbol{\theta}) = \mathbb{E}_{Z \sim p}[h(Z, \boldsymbol{\theta})] = \mathbb{E}_{Z \sim p}[g(T(Z, \boldsymbol{\theta}))],$$

where $h(\mathbf{z}, \boldsymbol{\theta}) := g(T(\mathbf{z}, \boldsymbol{\theta}))$. This implies

$$\begin{aligned}\nabla E(\boldsymbol{\theta}) &= \mathbb{E}_{Z \sim p}[\nabla_2 h(Z, \boldsymbol{\theta})] \\ &= \mathbb{E}_{Z \sim p}[\partial_2 T(Z, \boldsymbol{\theta})^* \nabla g(T(Z, \boldsymbol{\theta}))].\end{aligned}$$

The path gradient estimator (a.k.a. reparametrization trick) gives an **unbiased estimator** of $\nabla E(\boldsymbol{\theta})$. It has however two key disadvantages. First, it assumes that g is **differentiable** (almost everywhere), which may not always be the case. Second, it assumes that g is well-defined on \mathbb{R}^M , not on \mathcal{U} , which could be problematic for some discrete loss functions, such as the zero-one loss function or ranking loss functions.

As an example of differentiable transform, in machine learning, we can sample Gaussian noise Z and make it go through a neural network with parameters \mathbf{w} to generate an image $X := T(Z, \mathbf{w})$. In statistics, many distributions are related to each other through differentiable transforms, as we recall below.

Example 12.3 (Some differentiable transforms in statistics). We give below a non-exhaustive list of differentiable transform examples.

- If $X \sim \text{Normal}(\mu, \sigma^2)$, then $\exp(X) \sim \text{Lognormal}(\mu, \sigma^2)$.
- If $U \sim \text{Uniform}(0, 1)$, then $-\log(U)/\lambda \sim \text{Exponential}(\lambda)$.
- If $X_1, \dots, X_N \sim \text{Exponential}(\lambda)$ (i.i.d.), then $\sum_{i=1}^N X_i \sim \text{Gamma}(N, \lambda)$.
- If $X_i \sim \text{Gamma}(\alpha_i, \theta)$ for $i \in [K]$, then $\left(\frac{X_1}{\sum_{i=1}^K X_i}, \dots, \frac{X_K}{\sum_{i=1}^K X_i} \right) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$.

12.4.3 Inverse transforms

The inverse transform method can be used for sampling from a probability distribution, given access to its associated **quantile function**. Recall that the cumulative distribution function (CDF) associated with a random variable Y is the function $F_Y: \mathbb{R} \rightarrow [0, 1]$ defined by

$$F_Y(y) := \mathbb{P}(Y \leq y).$$

The quantile function is then a function $Q_Y: [0, 1] \rightarrow \mathbb{R}$ such that $Q_Y(\pi) = y$ for $\pi = F_Y(y)$. Assuming F_Y is continuous and strictly increasing, we have that Q_Y is the **inverse CDF**,

$$Q_Y(\pi) = F_Y^{-1}(\pi).$$

In the general case of CDF functions that are not strictly increasing, the quantile function is usually defined as

$$Q_Y(\pi) := \inf\{y \in \mathbb{R}: \pi \leq F_Y(y)\}.$$

Given access to the quantile function $Q_Y(\pi)$ associated with a distribution p , inverse transform sampling allows us to sample from p by first drawing a sample from the **uniform distribution** and then making this sample go through the quantile function.

Proposition 12.6 (Inverse transform sampling). Suppose $Y \sim p$, where p is a distribution with quantile function Q_Y . If $U \sim \text{Uniform}(0, 1)$, then $Q_Y(U) \sim p$.

Proof. If $\pi \leq F_Y(t)$, then by definition of Q_Y , $Q_Y(\pi) \leq t$. If $\pi \geq F_Y(t)$, then by definition of Q_Y , $F_Y(Q_Y(\pi)) \geq \pi$, so $F_Y(Q_Y(\pi)) \geq F_Y(t)$ and since a CDF is always non-decreasing, $Q_Y(\pi) \geq t$. Hence, we have, $Q_Y(\pi) \leq t \iff \pi \leq F_Y(t)$, so

$$\begin{aligned} \mathbb{P}(Q_Y(U) \leq t) &= \mathbb{P}(U \leq F_Y(t)) \\ &= F_Y(t). \end{aligned}$$

The CDFs of $Q_Y(U)$ and Y coincide, hence they have the same distribution. \square

If the quantile function is differentiable, we can therefore use it as a **transformation** within the **reparametrization trick**. Indeed, if $Y \sim p_{\theta}$, where p_{θ} is a distribution with parameter θ and quantile function $Q_Y(\pi, \theta)$, then we have

$$E(\theta) = \mathbb{E}_{Y \sim p_{\theta}}[g(Y)] = \mathbb{E}_{\pi \sim \text{Uniform}(0,1)}[g(Q_Y(\pi, \theta))]$$

and therefore, by the reparametrization trick (Proposition 12.5),

$$\nabla E(\theta) = \mathbb{E}_{\pi \sim \text{Uniform}(0,1)}[\partial_2 Q_Y(\pi, \theta)^* \nabla g(Q_Y(\pi, \theta))].$$

Example 12.4 (Examples of quantile functions). If

$Y \sim \text{Exponential}(\lambda)$, the CDF of Y is $\pi = F_Y(y) = 1 - \exp(-\lambda y)$ for $y \geq 0$ and therefore the quantile function is $Q_Y(\pi, \lambda) = -\frac{\log(1-\pi)}{\lambda}$. If $Y \sim \text{Normal}(\mu, \sigma^2)$, the CDF is $F_Y(y) = \frac{1}{2} \left[1 + \text{erf} \left(\frac{y-\mu}{\sigma\sqrt{2}} \right) \right]$ and the quantile function is $Q_Y(\pi, \theta) = \mu + \sigma\sqrt{2} \cdot \text{erf}^{-1}(2\pi - 1)$, where $\theta = (\mu, \sigma)$. This therefore defines an alternative transformation to the location-scale transformation in Eq. (12.4).

Note that, in the above example, the error function erf and its inverse do not enjoy analytical expressions but autodiff packages usually provide numerical routines to compute them and differentiate through them. Nonetheless, one caveat of the inverse transform is that it indeed requires access to (approximations of) the quantile function and its derivatives, which may be difficult for complicated distributions.

12.4.4 Pushforward operators

Pushforward distributions

We saw so far that the reparametrization trick is based on using a change of variables in order to differentiate an expectation w.r.t. the parameters of the distribution. In this section, we further formalize that approach using pushforward distributions.

Definition 12.1 (Pushforward distribution). Suppose $Z \sim p$, where p is a distribution over \mathcal{Z} . Given a continuous map $T: \mathcal{Z} \rightarrow \mathcal{U}$, the pushforward distribution of p through T is the distribution q according to which $U := T(Z) \in \mathcal{U}$ is distributed, i.e., $U \sim q$.

Although not explicit in the above, the transformation T can depend on some learnable parameters, for example if T is a neural network. Intuitively, the pushforward distribution is obtained by moving the position of all the points in the support of p . We give a few examples below.

- Inverse transform sampling studied in Section 12.4.3 can be seen as performing the pushforward of the uniform distribution through $T = Q$, where Q is the quantile function.

- The Gumbel trick studied in Section 14.5 can be seen as a the push-forward of Gumbel noise through $T = \operatorname{argmax}$ (a discontinuous function).
- Gumbel noise can itself be obtained by pushing forward the uniform distribution through $T = -\log(-\log(\cdot))$ (Remark 14.3).
- In a generative modeling setting, as we mentioned previously, we use the pushforward of Gaussian noise through a parametrized transformation $X = T(Z, \mathbf{w})$ called a generator, typically a neural network.
- It is possible to define distributions over (sparse) probability vectors by sampling then projecting (Farinhas *et al.*, 2021).

A crucial aspect of the pushforward distribution q is that it can be **implicitly** defined, meaning that we do not necessarily need to know the explicit form of the associated PDF. In fact, it is easy to **sample** from q , provided that it is easy to sample from p :

$$U \sim q \iff Z \sim p, U := T(Z).$$

Hence the usefulness of the pushforward distribution in **generative modeling**. Furthermore, if p has associated PDF p_Z , we can compute the expectation of a function f according to q as

$$\mathbb{E}_{U \sim q}[f(U)] = \mathbb{E}_{Z \sim p}[f(T(Z))] = \int_{\mathcal{Z}} f(T(\mathbf{z}))p_Z(\mathbf{z})d\mathbf{z},$$

even though we do not know the explicit form of the PDF of q .

Pushforward measures

More generally, we can define the notion of pushforward, in the language of measures. Denote $\mathcal{M}(\mathcal{Z})$ the set of measures on a set \mathcal{Z} . A **measure** $\alpha \in \mathcal{M}(\mathcal{Z})$, that has a density $d\alpha(\mathbf{z}) := p_Z(\mathbf{z})d\mathbf{z}$, can be integrated against a function f as

$$\int_{\mathcal{Z}} f(\mathbf{z})d\alpha(\mathbf{z}) = \int_{\mathcal{Z}} f(\mathbf{z})p_Z(\mathbf{z})d\mathbf{z}.$$

A measure α is called a probability measure if it is positive and satisfies $\alpha(\mathcal{Z}) = \int_{\mathcal{Z}} d\alpha(\mathbf{z}) = \int_{\mathcal{Z}} p_Z(\mathbf{z}) d\mathbf{z} = 1$. See Peyré and Cuturi (2019, Chapter 2) for a concise introduction.

Definition 12.2 (Pushforward operator and measure). Given a continuous map $T: \mathcal{Z} \rightarrow \mathcal{U}$ and some measure $\alpha \in \mathcal{M}(\mathcal{Z})$, the pushforward measure $\beta = T_{\#}\alpha \in \mathcal{M}(\mathcal{U})$ is such that for all continuous functions $f \in \mathcal{C}(\mathcal{U})$

$$\int_{\mathcal{U}} f(\mathbf{u}) d\beta(\mathbf{u}) = \int_{\mathcal{Z}} f(T(\mathbf{z})) d\alpha(\mathbf{z}).$$

Equivalently, for any measurable set $\mathcal{A} \subset \mathcal{U}$, we have

$$\beta(\mathcal{A}) = \alpha(\{\mathbf{z} \in \mathcal{Z} : T(\mathbf{z}) \in \mathcal{A}\}) = \alpha(T^{-1}(\mathcal{A})),$$

where $T^{-1}(\mathcal{A}) = \{\mathbf{z} \in \mathcal{Z} : T(\mathbf{z}) \in \mathcal{A}\}$.

Importantly, the pushforward operator preserves positivity and mass, therefore if α is a probability measure, then so is $T_{\#}\alpha$. The pushforward of a probability measure therefore defines a pushforward distribution (since a distribution can be parametrized by a probability measure).

12.4.5 Change-of-variables theorem

We saw that a pushforward distribution associated with a variable U is implicitly defined through a transform $U := T(Z)$ and can be easily sampled from as long as it is easy to sample Z . However, in some applications (e.g., density estimation), we may want to know the PDF associated with U . Assuming the transform T is invertible, we have $Z = T^{-1}(U)$ and therefore for $\mathcal{A} \subseteq \mathcal{U}$, we have

$$\mathbb{P}(U \in \mathcal{A}) = \mathbb{P}(Z \in T^{-1}(\mathcal{A})) = \int_{T^{-1}(\mathcal{A})} p_Z(\mathbf{z}) d\mathbf{z}.$$

Using the **change-of-variables theorem** from multivariate calculus, assuming T^{-1} is available, we can give an explicit formula for the PDF of the pushforward distribution, see e.g. (Schwartz, 1954; Taylor, 2002).

Proposition 12.7 (PDF of the pushforward distribution). Suppose $Z \sim p$, where p is a distribution over \mathcal{Z} , with PDF p_Z . Given a **diffeomorphism** $T: \mathcal{Z} \rightarrow \mathcal{U}$ (i.e., an invertible and differentiable map), the pushforward distribution of p through T is the distribution q such that $U := T(Z) \sim q$ and its PDF is

$$q_U(\mathbf{u}) = |\det(\partial T^{-1}(\mathbf{u}))| p_Z(T^{-1}(\mathbf{u})),$$

where $\partial T^{-1}(\mathbf{u})$ is the Jacobian of $T^{-1}: \mathcal{U} \rightarrow \mathcal{Z}$.

Using this formula, we obtain

$$\begin{aligned} \mathbb{P}(U \in \mathcal{A}) &= \int_{\mathcal{A}} p_U(\mathbf{u}) d\mathbf{u} \\ &= \int_{\mathcal{A}} |\det(\partial T^{-1}(\mathbf{u}))| p_Z(T^{-1}(\mathbf{u})) d\mathbf{u}. \end{aligned}$$

Using the inverse function theorem (Theorem 11.5), we then have

$$\partial T^{-1}(\mathbf{u}) = (\partial T(T^{-1}(\mathbf{u})))^{-1},$$

under the assumption that $T(\mathbf{z})$ is continuously differentiable and has invertible Jacobian $\partial T(\mathbf{z})$. **Normalizing flows** are parametrized transformations T designed such that T^{-1} and its Jacobian ∂T^{-1} are easy to compute; see e.g. Kobyzev *et al.* (2019) and Papamakarios *et al.* (2021) for a review.

12.5 Stochastic programs

A stochastic program is a program that involves some form of randomness. In a stochastic program, the final output, as well as intermediate variables, may therefore be random variables. In other words, a stochastic program induces a probability distribution over program outputs, as well as over execution trajectories.

12.5.1 Stochastic computation graphs

A stochastic program can be represented by a stochastic computation graph as originally introduced by Schulman *et al.* (2015). Departing from that work, our exposition explicitly supports two types of intermediate

operations: sampling from a **conditional distribution** or evaluating a **function**. These operations can produce either **deterministic** variables or **random** variables.

Function and distribution nodes

Formally, we define a stochastic computation graph as a directed acyclic graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \mathcal{V}_f \cup \mathcal{V}_p$, \mathcal{V}_f is the set of function nodes and \mathcal{V}_p is the set of distribution nodes. Similarly to computation graphs reviewed in Section 4.1.3, we number the nodes as $\mathcal{V} = \{0, 1, \dots, K\}$. Node 0 corresponds to the input $s_0 \in \mathcal{S}_0$, which we assume to be deterministic. It is the variable with respect to which we wish to differentiate. Node K corresponds to the program output $S_K \in \mathcal{S}_K$, which we assume to be a random variable. A node $k \in \{1, \dots, K\}$ can either be a **function node** $k \in \mathcal{V}_f$ with an associated function f_k or a **distribution node** $k \in \mathcal{V}_p$, with associated conditional distribution p_k . A stochastic program has at least one distribution node, the source of randomness. Otherwise, it is a deterministic program. As for computation graphs, the set of edges \mathcal{E} is used to represent dependencies between nodes. We denote the parents of node k by $\text{pa}(k)$.

Deterministic and random variables

We distinguish between two types of intermediate variables: **deterministic** variables s_k and **random** variables S_k . Therefore, a distribution p_k or a function f_k may receive both types of variables as **conditioning** or **input**. It is then convenient to split $\text{pa}(k)$ as $\text{pa}(k) = \text{determ}(k) \cup \text{random}(k)$, where we defined the **deterministic parents** $\text{determ}(k) := \{i_1, \dots, i_{p_k}\}$ and the **random parents** $\text{random}(k) := \{j_1, \dots, j_{q_k}\}$. Therefore, $s_{i_1}, \dots, s_{i_{p_k}}$ are the deterministic parent variables and $S_{j_1}, \dots, S_{j_{q_k}}$ are the random parent variables, of node k .

Executing a stochastic program

We assume that nodes $0, 1, \dots, K$ are in topological order (if this is not the case, we need to perform a topological sort). Given parent

variables $\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}}$ and $S_{j_1}, \dots, S_{j_{q_k}}$, a node $k \in \{1, \dots, K\}$ produces an output as follows.

- If $k \in \mathcal{V}_p$ (distribution node), the output is

$$\begin{aligned} S_k &\sim p_k(\cdot \mid \mathbf{s}_{\text{determin}(k)}, S_{\text{random}(k)}) \\ \iff S_k &\sim p_k(\cdot \mid \mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}}, S_{j_1}, \dots, S_{j_{q_k}}) \end{aligned}$$

Note that technically p_k is the distribution of S_k conditioned on its parents, not the distribution of S_k . Therefore, we should in principle write $S_k \mid \mathbf{s}_{\text{determin}(k)}, S_{\text{random}(k)} \sim p_k(\cdot \mid \mathbf{s}_{\text{determin}(k)}, S_{\text{random}(k)})$. We avoid this notation for conciseness and for symmetry with function nodes.

Contrary to a function node, a distribution node can have no parents. That is, if $k \in \mathcal{V}_p$, it is possible that $\text{pa}(k) = \emptyset$. A good example would be a parameter-free noise distribution.

- If $k \in \mathcal{V}_f$ (function node), the output is in general

$$\begin{aligned} S_k &:= f_k(\mathbf{s}_{\text{determin}(k)}, S_{\text{random}(k)}) \\ &:= f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}}, S_{j_1}, \dots, S_{j_{q_k}}) \end{aligned}$$

and in the special case $q_k = |\text{random}(k)| = 0$, the output is

$$\begin{aligned} \mathbf{s}_k &:= f_k(\mathbf{s}_{\text{determin}(k)}) \\ &:= f_k(\mathbf{s}_{i_1}, \dots, \mathbf{s}_{i_{p_k}}). \end{aligned}$$

Unless the associated conditional distribution p_k is a delta distribution, that puts all the probability mass on a single point, the output of a distribution node $k \in \mathcal{V}_p$ is necessarily a random variable $S_k \in \mathcal{S}_k$. For function nodes $k \in \mathcal{V}_f$, the output of the function f_k is a random variable $S_k \in \mathcal{S}_k$ if at least one of the parents of k produces a random variable. Otherwise, if all parents of k produce deterministic variables, the output of f_k is a deterministic variable $\mathbf{s}_k \in \mathcal{S}_k$.

The entire procedure is summarized in Algorithm 12.1. We emphasize that $S_K = f(\mathbf{s}_0) \in \mathcal{S}_K$ is a random variable. Therefore, a stochastic program (implicitly) induces a distribution over S_K , and also over intermediate random variables S_k . Executing the stochastic program allows us to draw samples from that distribution.

Algorithm 12.1 Executing a stochastic program

Nodes: $1, \dots, K$ in topological order, where node k is either a function f_k or a conditional distribution p_k

Input: input $s_0 \in \mathcal{S}_0$

- 1: **for** $k := 1, \dots, K$ **do**
- 2: Retrieve $\text{pa}(k) = \text{determ}(k) \cup \text{random}(k)$
- 3: **if** $k \in \mathcal{V}_p$ **then** ▷ Distribution node
- 4: $S_k \sim p_k(\cdot | s_{\text{determ}(k)}, S_{\text{random}(k)})$
- 5: **else if** $k \in \mathcal{V}_f$ **then** ▷ Function node
- 6: **if** $|\text{random}(k)| \neq 0$ **then**
- 7: $S_k := f_k(s_{\text{determ}(k)}, S_{\text{random}(k)})$ ▷ Output is a R.V.
- 8: **else if** $|\text{random}(k)| = 0$ **then**
- 9: $s_k := f_k(s_{\text{determ}(k)})$ ▷ Output is deterministic
- 10: **Output:** $f(s_0) := S_K \in \mathcal{S}_K$

Special cases

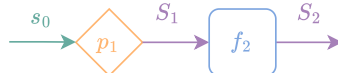
If all nodes are function nodes, we recover computation graphs, reviewed in Section 4.1.3. If all nodes are distribution nodes, we recover Bayesian networks, reviewed in Section 10.5.

12.5.2 Examples

We now present several examples that illustrate our formalism. We use the legend below in the following illustrations.

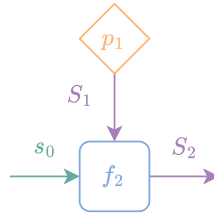


- Example 1 (SFE estimator):



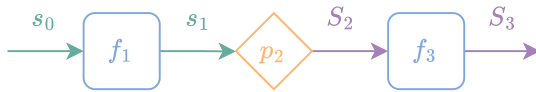
$$\begin{aligned}
S_1 &\sim p_1(\cdot \mid \mathbf{s}_0) \\
S_2 &:= f_2(S_1) \\
E(\mathbf{s}_0) &:= \mathbb{E}[S_2] \\
\nabla E(\mathbf{s}_0) &= \mathbb{E}_{S_1}[f_2(S_1) \nabla_{\mathbf{s}_0} \log p_1(S_1 \mid \mathbf{s}_0)]
\end{aligned}$$

- Example 2 (Pathwise estimator):



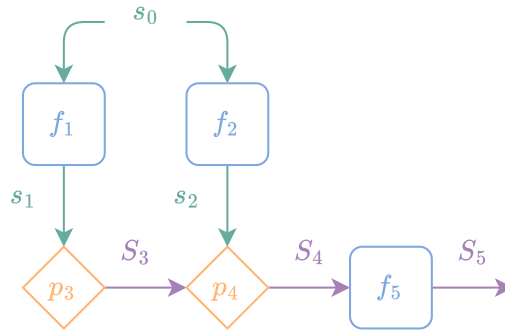
$$\begin{aligned}
S_1 &\sim p_1 \\
S_2 &:= f_2(S_1, \mathbf{s}_0) \\
E(\mathbf{s}_0) &:= \mathbb{E}[S_2] \\
\nabla E(\mathbf{s}_0) &= \mathbb{E}_{S_1} [\nabla_{\mathbf{s}_0} f_2(S_1, \mathbf{s}_0)]
\end{aligned}$$

- Example 3 (SFE estimator + chain rule):



$$\begin{aligned}
\mathbf{s}_1 &:= f_1(\mathbf{s}_0) \\
S_2 &\sim p_2(\cdot \mid \mathbf{s}_1) \\
S_3 &:= f_3(S_2) \\
E(\mathbf{s}_0) &:= \mathbb{E}[S_3] \\
\nabla E(\mathbf{s}_0) &= \partial f(\mathbf{s}_0)^* \mathbb{E}_{S_2}[f_3(S_2) \nabla_{\mathbf{s}_1} \log p_2(S_2 \mid \mathbf{s}_1)]
\end{aligned}$$

- Example 4:



$$\mathbf{s}_1 := f_1(\mathbf{s}_0)$$

$$\mathbf{s}_2 := f_2(\mathbf{s}_0)$$

$$S_3 \sim p_3(\cdot \mid \mathbf{s}_1)$$

$$S_4 \sim p_4(\cdot \mid \mathbf{s}_2, S_3)$$

$$S_5 := f_5(S_4)$$

$$E(\mathbf{s}_0) := \mathbb{E}[S_5] = \mathbb{E}_{S_3} [\mathbb{E}_{S_4} [f_5(S_4)]]$$

$$\begin{aligned} \nabla E(\mathbf{s}_0) = & \mathbb{E}_{S_3} [\partial f_1(\mathbf{s}_0)^* \nabla_{\mathbf{s}_1} \log p(S_3 \mid \mathbf{s}_1) \mathbb{E}_{S_4} [f_5(S_4)]] \\ & + \mathbb{E}_{S_3} [\mathbb{E}_{S_4} [\partial f_2(\mathbf{s}_0)^* \nabla_{\mathbf{s}_2} \log p_4(S_4 \mid \mathbf{s}_2, S_3) f_5(S_4)]] \end{aligned}$$

As can be seen, the gradient expressions can quickly become quite complicated, demonstrating the merits of automatic differentiation in stochastic computation graphs.

12.5.3 Unbiased gradient estimators

The output of a stochastic program is a random variable

$$S_K := f(\mathbf{s}_0).$$

It implicitly defines a probability distribution $p(\cdot \mid \mathbf{s}_0)$ such that $S_K \sim p(\cdot \mid \mathbf{s}_0)$. Executing the stochastic program once gives us an i.i.d. sample from $p(\cdot \mid \mathbf{s}_0)$.

Since derivatives are defined for deterministic variables, we need a way to convert a random variable to a deterministic variable. One way to do so is to consider the expected value (another way would be the mode)

$$E(\mathbf{s}_0) := \mathbb{E}[S_K] = \mathbb{E}[f(\mathbf{s}_0)] \in \text{conv}(\mathcal{S}_K),$$

where the expectation is over $S_K \sim p(\cdot | \mathbf{s}_0)$ or equivalently over the intermediate random variables S_k

$$S_k \sim p_k(\cdot | \mathbf{s}_{\text{determin}(k)}, S_{\text{random}(k)}),$$

for $k \in \mathcal{V}_p$ (the distribution nodes). We then wish to compute the gradient or more generally the Jacobian of $E(\mathbf{s}_0)$.

If all nodes in the stochastic computation graph are function nodes, we can estimate the gradient of $E(\mathbf{s}_0)$ using the pathwise estimator a.k.a. reparametrization trick (Section 12.4). This is the approach taken by Kingma and Welling (2013) and Rezende *et al.* (2014).

If all nodes in the stochastic computation graph are distribution nodes, we can use the SFE estimator (Section 12.3). Schulman *et al.* (2015) propose a surrogate loss so that using autodiff on that loss produces an unbiased gradient of the expectation, using the SFE estimator. Foerster *et al.* (2018) extend the approach to support high-order differentiation. Krieken *et al.* (2021) further extend the approach by supporting different estimators per node, as well as control variates.

Converting distribution nodes into function nodes and vice-versa

Our formalism uses two types of nodes: distribution nodes with associated conditional distribution p_k and function nodes with associated function f_k . It is often possible to convert between node types.

Converting a distribution node into a function node is exactly the reparametrization trick studied in Section 12.4. We can use transformations such as the location-scale transform or the inverse transform.

Converting a function node into a distribution node can be done using the change-of-variables theorem, studied in Section 12.4.5, on a pushforward distribution.

Because the pathwise estimator has lower variance than SFE, this is the method of choice when the f_k functions are available. The conversion from distribution node to function node and vice-versa is illustrated in Fig. 12.1.

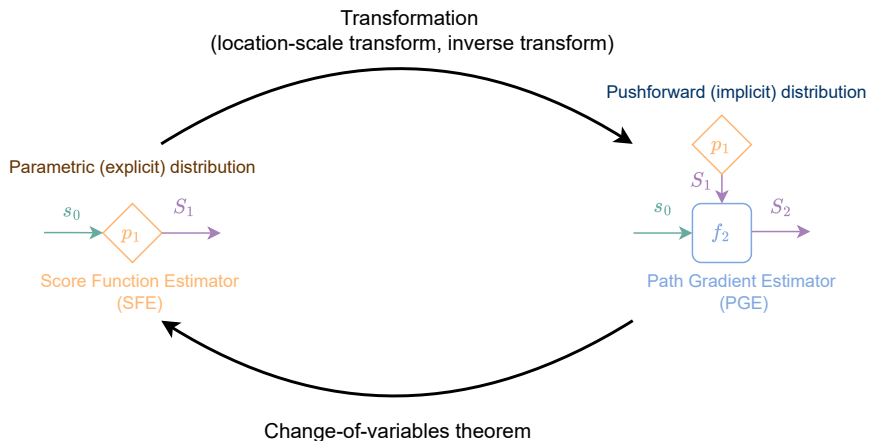


Figure 12.1: It is sometimes possible to convert a distribution node to a function node and vice-versa using a suitable transformation.

12.5.4 Local vs. global expectations

A stochastic computation graph can be seen as a **stochastic process**, a collection of random variables S_k , indexed by k , the position in the topological order. However, random variables are incompatible with autodiff. Replacing random variables by their expectation can be seen as a way to make them compatible with autodiff. Two strategies are then possible.

As we saw in the previous section, a strategy is to consider the expectation of the last output S_K . This strategy corresponds to a **global smoothing**. The two major advantages are that i) we do not need to assume that f_{k+1} is well-defined on $\text{conv}(\mathcal{S}_k)$ and ii) this induces a probability distribution over program executions. This is for instance useful to compute the variance of the program. The gradient of the program's expected value can be estimated by the reparametrization trick or by the SFE, depending on the type of nodes used.

A second strategy is to replace an intermediate random variable $S_k \in \mathcal{S}_k$, for $k \in \{1, \dots, K\}$, by its expectation $\mathbb{E}[S_k] \in \text{conv}(\mathcal{S}_k)$. This strategy corresponds to a **local smoothing**. A potential drawback of this approach is that $\mathbb{E}[S_k]$ belongs to $\text{conv}(\mathcal{S}_k)$, the convex hull of \mathcal{S}_k .

Therefore, the function f_{k+1} in which $\mathbb{E}[S_k]$ is fed must be well-defined on $\text{conv}(\mathcal{S}_k)$, which may not always be the case. In the case of control flows, another disadvantage is computational. We saw in Section 5.6 and Section 5.7 that using a soft comparison operator within a conditional statement induces a distribution on a binary or categorical random variable, corresponding to the branch to be selected. A conditional statement can then be locally smoothed out by replacing the random variable by its expectation i.e., a **convex combination** of all the branches. This means that, unless the distribution has sparse support, all branches must be evaluated.

12.6 Differential equations

12.6.1 Parameterized differential equations

From residual networks to neural ODEs

Starting from $\mathbf{s}_0 := \mathbf{x}$, residual networks, reviewed in Section 4.6, iterate for $k \in \{1, \dots, K\}$

$$\mathbf{s}_k := \mathbf{s}_{k-1} + h_k(\mathbf{s}_{k-1}, \mathbf{w}_k).$$

A residual network can be seen as parameterizing incremental discrete-time input changes (hence the name “residual”)

$$\mathbf{s}_k - \mathbf{s}_{k-1} = h_k(\mathbf{s}_k, \mathbf{w}_k).$$

Chen *et al.* (2018) proposed to parameterize continuous-time (instantaneous) changes instead. They considered the evolution $\mathbf{s}(t)$ of the inputs in continuous time driven by a function $h(t, \mathbf{s}, \mathbf{w})$ parameterized by \mathbf{w} , starting from \mathbf{x} . Formally, the evolution $\mathbf{s}(t)$ is the solution of the **ordinary differential equation** (ODE)

$$\begin{aligned} \mathbf{s}(0) &= \mathbf{x} \\ \mathbf{s}'(t) &= h(t, \mathbf{s}(t), \mathbf{w}) \quad t \in [0, T] \end{aligned} \tag{12.5}$$

Here, $\mathbf{s}'(t)$ is the vector of derivatives of \mathbf{s} as defined in Remark 2.4, and T denotes a final time for the trajectory. The output of such a **neural ODE** (Chen *et al.*, 2018) is then $f(\mathbf{x}, \mathbf{w}) := \mathbf{s}(T)$. Alternatively, the

output can be seen as the solution of an **integration** problem

$$f(\mathbf{x}, \mathbf{w}) = \mathbf{s}(T) = \mathbf{x} + \int_0^T h(t, \mathbf{s}(t), \mathbf{w}) dt. \quad (12.6)$$

Differential equations like Eq. (12.5) arise in many contexts beyond neural ODEs, ranging from modeling physical systems to pandemics (Braun and Golubitsky, 1983). Moreover, the differential equation presented in Eq. (12.5) is just an example of an ordinary differential equation, while controlled differential equations or stochastic differential equations can also be considered.

Existence of a solution

First and foremost, the question is whether $\mathbf{s}(t)$ is well-defined. Fortunately, the answer is positive under mild conditions, as shown by Picard-Lindelöf's theorem recalled below (Butcher, 2016, Theorem 16).

Theorem 12.1 (Existence and uniqueness of ODE solutions). If $h : [0, T] \times \mathcal{S} \rightarrow \mathcal{S}$ is continuous in its first variable and Lipschitz-continuous in its second variable, then there exists a unique differentiable map $\mathbf{s} : [0, T] \rightarrow \mathcal{S}$ satisfying

$$\begin{aligned} \mathbf{s}(0) &= \mathbf{s}_0 \\ \mathbf{s}'(t) &= h(t, \mathbf{s}(t)) \quad t \in [0, T], \end{aligned}$$

for some given $\mathbf{s}_0 \in \mathcal{S}$.

For time-independent linear functions $h(t, \mathbf{s}) = \mathbf{A}\mathbf{s}$, the integral in Eq. (12.6) can be computed in closed form as

$$\mathbf{s}_t = \exp(t\mathbf{A})(\mathbf{s}_0),$$

where $\exp(\mathbf{A})$ is the matrix exponential. Hence, the output $\mathbf{s}(T)$ can be expressed as a simple function of the parameters (\mathbf{A} in this case). However, generally, we do not have access to such analytical solutions, and, just as for solving optimization problems in Chapter 11, we need to resort to some iterative algorithms.

Integration methods

To numerically solve an ODE, we can use **integration methods**, whose goal is to build a sequence \mathbf{s}_k that approximates the solution $\mathbf{s}(t)$ at times t_k . The simplest integration method is the **explicit Euler method**, that approximates the solutions between times t_{k-1} and t_k as

$$\begin{aligned}\mathbf{s}(t_{k-1}) - \mathbf{s}(t_k) &= \int_{t_{k-1}}^{t_k} h(t, \mathbf{s}(t), \mathbf{w}) dt \\ &\approx \delta_k h(t_{k-1}, \mathbf{s}(t_{k-1}), \mathbf{w}),\end{aligned}$$

for a time-step

$$\delta_k := t_k - t_{k-1}.$$

The resulting integration scheme consists in computing starting from $\mathbf{s}_0 = \mathbf{x}$, for $k \in \{1, \dots, K\}$,

$$\mathbf{s}_k := \mathbf{s}_{k-1} + \delta_k h(t_{k-1}, \mathbf{s}_{k-1}, \mathbf{w}).$$

Assimilating $\delta_k h(t_{k-1}, \mathbf{s}_{k-1}, \mathbf{w})$ with $h_k(\mathbf{s}_{k-1}, \mathbf{w}_k)$, we find that residual networks are essentially the discretization of a neural ODE by an explicit Euler method; more precisely, a non-autonomous neural ODEs, see e.g. (Davis *et al.*, 2020).

Euler's forward method is only one integration method among many. To cite a few, there are implicit Euler methods, semi-implicit methods, Runge-Kutta methods, linear multistep methods, etc. See, e.g., Gautschi (2011) for a detailed review. The quality of an integration method is measured by its consistency and its stability (Gautschi, 2011). These concepts naturally influence the development of evaluation and differentiation techniques for ODEs. We briefly summarize them below.

Given a fixed time interval $\delta_k = \delta$ and $K = \lceil T/\delta \rceil$ points, an integration method is **consistent of order k** if $\|\mathbf{s}_k - \mathbf{s}(k\delta)\| = O(\delta^k)$ as $\delta \rightarrow 0$ and therefore $k \rightarrow +\infty$. The higher the order k , the fewer points we need to reach an approximation error ε on the points considered. The term $\|\mathbf{s}_k - \mathbf{s}(k\delta)\| = O(\delta^k)$ is reminiscent of the error encountered in finite differences (Chapter 7) and is called the **truncation error**.

The (absolute) **stability** of a method is defined by the set of time-steps such that the integration method can integrate $s'(t) = \lambda s(t)$ for some $\lambda \in \mathbb{C}$ without blowing up as $t \rightarrow +\infty$.

12.6.2 Continuous adjoint method

Since different parameters \mathbf{w} induce different trajectories associated to $h(t, \mathbf{s}, \mathbf{w})$ in Eq. (12.5), we may want to select one of these trajectories by minimizing some criterion. For example, we may consider selecting $\mathbf{w} \in \mathcal{W}$ by minimizing a loss L on the final point of the trajectory,

$$\min_{\mathbf{w} \in \mathcal{W}} L(f(\mathbf{x}, \mathbf{w}), \mathbf{y}), \quad (12.7)$$

where

$$f(\mathbf{x}, \mathbf{w}) := \mathbf{s}(T) = \mathbf{x} + \int_0^T h(t, \mathbf{s}(t), \mathbf{w}) dt.$$

To solve such problems, we need to access gradients of ℓ composed with f through VJPs of the solution of the ODE. The VJPs can actually be characterized as solutions of an ODE themselves thanks to the **continuous time adjoint method** (Pontryagin, 1985), presented below, and whose proof is postponed to Section 12.6.6.

Proposition 12.8 (Continuous-time adjoint method). Consider a function $h : [0, T] \times \mathcal{S} \times \mathcal{W} \rightarrow \mathcal{S}$, continuous in its first variable, Lipschitz-continuous and continuously differentiable in its second variable. Assume that $\partial_3 h(t, \mathbf{s}, \mathbf{w})$ exists for any $t, \mathbf{s}, \mathbf{w}$, and is also continuous in its first variable, Lipschitz-continuous in its second variable. Denote $\mathbf{s} : \mathcal{S} \rightarrow \mathcal{S}$ the solution of the ODE

$$\begin{aligned} \mathbf{s}(0) &= \mathbf{x} \\ \mathbf{s}'(t) &= h(t, \mathbf{s}(t), \mathbf{w}) \quad t \in [0, T], \end{aligned}$$

and $f(\mathbf{x}, \mathbf{w}) = \mathbf{s}(T)$ the final state of the ODE at time T .

Then, the function f is differentiable, and for an output direction $\mathbf{u} \in \mathcal{S}$, its VJP along \mathbf{u} is given by

$$\partial f(\mathbf{x}, \mathbf{w})^* \mathbf{u} = (\mathbf{r}(0), \mathbf{g})$$

for

$$\mathbf{g} = \int_0^T \partial_3 h(t, \mathbf{s}(t), \mathbf{w})^* \mathbf{r}(t) dt$$

and for \mathbf{r} solving the **adjoint** (backward) ODE

$$\begin{aligned}\mathbf{r}'(t) &= -\partial_2 h(t, \mathbf{s}(t), \mathbf{w})^* \mathbf{r}(t) \\ \mathbf{r}(T) &= \mathbf{u}.\end{aligned}$$

In particular, the gradient $\nabla(L \circ f)(\mathbf{x}, \mathbf{w})$ for $L : \mathcal{S} \rightarrow \mathbb{R}$ a differentiable loss is obtained by solving the adjoint ODE with $\mathbf{r}(T) = \nabla L(\mathbf{s}(T))$.

Example 12.5 (Fitting data through the solution of an ODE). As an illustrative example, we can consider optimizing the parameters of an ODE to fit some data points. Namely, we may seek a continuous time solution $z(t; \mathbf{w})$ of a modified Lotka Volterra ODE

$$\mathbf{z}'(t; \mathbf{w}) = \begin{pmatrix} \alpha z_1(t; \mathbf{w}) - \beta z_1(t; \mathbf{w}) z_2(t; \mathbf{w}) \\ -\gamma z_2(t; \mathbf{w}) + \delta z_1(t; \mathbf{w}) z_2(t; \mathbf{w}) \end{pmatrix} + \mathbf{c},$$

for $\mathbf{w} = (\alpha, \beta, \gamma, \delta, \mathbf{c})$, that fits some observations $\mathbf{z}_1, \dots, \mathbf{z}_T$. The optimization problem consists then of

$$\min_{\mathbf{w}} \sum_{\tau=1}^T (z(t_j; \mathbf{w}) - z_j)^2,$$

and requires backpropagating through the solution $\mathbf{z}(\cdot; \mathbf{w})$ of the ODE w.r.t. to its candidate parameters \mathbf{w} . Fig. 12.2 illustrates such a problem with varying candidate parameters

12.6.3 Gradients via the continuous adjoint method

Proposition 12.8 gives a formal definition of the gradient. However, just as computing the mapping $f(\mathbf{x}, \mathbf{w})$ itself, computing its VJP or the gradient of $L \circ f$ requires solving an integration problem. Note that the integration of $\mathbf{r}(t)$ in Proposition 12.8 requires also values of $\mathbf{s}(t)$. Therefore, we need to integrate both $\mathbf{r}(t)$ and $\mathbf{s}(t)$. Such an approach is generally referred as **optimize-then-discretize** because we first formulate the gradient in continuous time (the “optimize part”) and then discretize the resulting ODE.

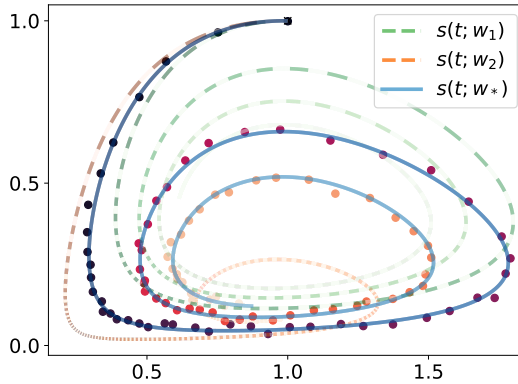


Figure 12.2: Finding the optimal parameters of an ODE to fit some observed data. The dots represent the trajectories of a dynamical system observed at regular times (time is represented here by a gradient color, the lighter the color, the larger the time). Each line represents the solution of an ODE given by some hyperparameters \mathbf{w} . The objective is to find the hyperparameters of the ODE such that its solution fits the data points. Green and orange lines fail to do so while the blue line fits the data. To compute such parameters \mathbf{w} , we need to backpropagate through the solution of the ODE.

Simple discretization scheme

A first approach consists in defining a backward discretization scheme that can approximate $\mathbf{s}(t)$ backward in time. Namely, by defining $\boldsymbol{\sigma}(t) = \mathbf{s}(T - t)$, $\boldsymbol{\rho}(t) = \mathbf{r}(T - t)$, and $\boldsymbol{\gamma}(t) = \int_t^T \partial_3 h(\tau, \mathbf{s}(\tau), \mathbf{w})^* \mathbf{r}(\tau) d\tau$, the derivative of $L \circ f$ is given by $(\boldsymbol{\rho}(T), \boldsymbol{\gamma}(T))$. The functions $\boldsymbol{\sigma}, \boldsymbol{\rho}, \boldsymbol{\gamma}$ are solutions of a standard ODE

$$\begin{aligned} \boldsymbol{\sigma}(0) &= \mathbf{s}(T), & \boldsymbol{\sigma}'(t) &= -h(T - t, \boldsymbol{\sigma}(t), \mathbf{w}), \\ \boldsymbol{\rho}(0) &= \nabla L(\mathbf{s}(T)), & \boldsymbol{\rho}'(t) &= \partial_2 h(T - t, \boldsymbol{\sigma}(t), \mathbf{w})^* \boldsymbol{\rho}(t), \\ \boldsymbol{\gamma}(0) &= 0, & \boldsymbol{\gamma}'(t) &= \partial_3 h(T - t, \boldsymbol{\sigma}(t), \mathbf{w})^* \boldsymbol{\rho}(t). \end{aligned}$$

The above ODE can then be solved by any integration method. Note, however, that it requires first computing $\mathbf{s}(T)$ and $\nabla L(\mathbf{s}(T))$ by an integration method. The overall computation of the gradient using an explicit Euler method to solve forward and backward ODEs is summarized in Algorithm 12.2.

Algorithm 12.2 naturally looks like the reverse mode of autodiff for a residual neural networks with **shared weights**. A striking difference

Algorithm 12.2 Gradient computation via continuous adjoint method with Euler explicit discretization

- 1: **Functions:** $h : [0, T] \times \mathcal{S} \times \mathcal{W} \rightarrow \mathbb{R}$, $L : \mathcal{S} \rightarrow \mathbb{R}$
 - 2: **Inputs:** input \mathbf{x} , parameters \mathbf{w} , number of discretization steps K .
 - 3: Set discretization step $\delta = T/K$, denote $h_k(\mathbf{s}, \mathbf{w}) = h(k\delta, \mathbf{s}, \mathbf{w})$.
 - 4: Set $\mathbf{s}_0 := \mathbf{x}$
 - 5: **for** $k := 1, \dots, K$ **do** ▷ Forward discretization
 - 6: Compute $\mathbf{s}_k := \mathbf{s}_{k-1} + \delta h_{k-1}(\mathbf{s}_{k-1}, \mathbf{w})$.
 - 7: Compute $\mathbf{u} := \nabla L(\mathbf{s}_K)$.
 - 8: Initialize $\mathbf{r}_K := \mathbf{u}$, $\hat{\mathbf{s}}_K = \mathbf{s}_K$, $\mathbf{g}_K = \mathbf{0}$
 - 9: **for** $k := K, \dots, 1$ **do** ▷ Backward discretization
 - 10: Compute $\hat{\mathbf{s}}_{k-1} := \hat{\mathbf{s}}_k - \delta h_k(\hat{\mathbf{s}}_k, \mathbf{w})$
 - 11: Compute $\mathbf{r}_{k-1} := \mathbf{r}_k + \delta \partial_2 h_k(\hat{\mathbf{s}}_k, \mathbf{w})^* \mathbf{r}_k$
 - 12: Compute $\mathbf{g}_{k-1} := \mathbf{g}_k + \delta \partial_3 h_k(\hat{\mathbf{s}}_k, \mathbf{w})^* \mathbf{r}_k$
 - 13: **Output:** $(\mathbf{r}_0, \mathbf{g}_0) \approx \nabla(L \circ f)(\mathbf{x}, \mathbf{w})$
-

is that the intermediate computations \mathbf{s}_k are not kept in memory and, instead, new variables $\hat{\mathbf{s}}_k$ are computed along the backward ODE. One may believe that by switching to continuous time, we solved the memory issues encountered in reverse-mode autodiff. Unfortunately, this comes at the cost of numerical stability. As we use a discretization scheme to recompute the intermediate states backward in time through $\hat{\mathbf{s}}_k$ in Algorithm 12.2, we accumulate some truncation errors.

To understand the issue here, consider applying Algorithm 12.2 repeatedly on the same parameters but using $\hat{\mathbf{s}}_0$ instead of $\mathbf{s}_0 = \mathbf{x}$ each time. In the continuous realm, $\sigma(T) = \mathbf{s}(0)$. But after discretization, $\hat{\mathbf{s}}_0 \approx \sigma(T)$ does not match \mathbf{s}_0 . Therefore, by applying Algorithm 12.2 with $\mathbf{s}_0 = \hat{\mathbf{s}}_0$, we would not get the same output even if in continuous time we naturally should have. This phenomenon is illustrated in Fig. 12.3. It intuitively shows why Algorithm 12.2 induces some noise in the estimation of the gradient.

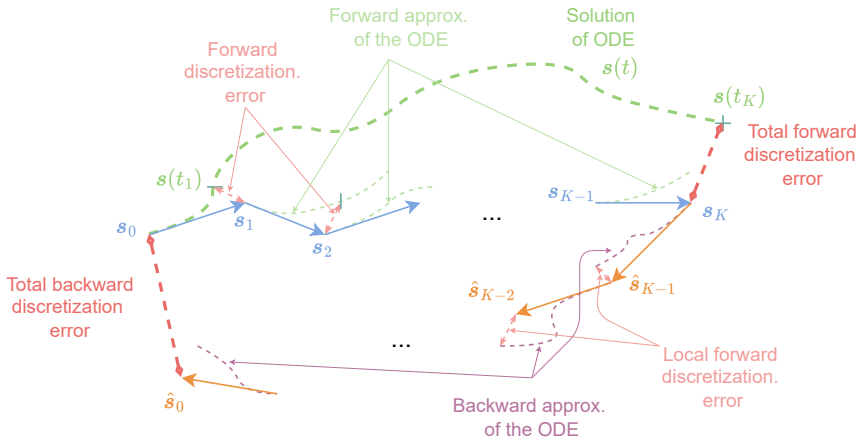


Figure 12.3: Forward and backward discretizations when using the continuous adjoint method.

Multiple shooting scheme

An alternative approach consists in integrating both the forward and backward ODEs jointly. Namely, we may solve an ODE with boundary values

$$\begin{aligned} s'(t) &= h(t, s(t), w), & s(0) &= x, \\ r'(t) &= -\partial_2 h(t, s(t), w)^* r(t), & r(T) &= \nabla L(s(T)) \\ g'(t) &= -\partial_3 h(t, s(t), w)^* r(t), & g(T) &= 0, \end{aligned}$$

by means of a multiple shooting method or a collocation method (Stoer *et al.*, 1980). This approach still requires $\nabla L(s(T))$ to be approximated first.

12.6.4 Gradients via reverse-mode on discretization

A simpler approach consists in replacing the objective in Eq. (12.7) by its version discretized using some numerical method, such as an Euler forward discretization scheme. That is, we seek to solve

$$\min_{w \in \mathcal{W}} L(s_K) \quad \text{where} \quad s_k = s_{k-1} + \delta h(k\delta, s_{k-1}, w) \quad k \in \{1, \dots, K\},$$

with $\mathbf{s}_0 = \mathbf{0}$ and δ some discretization step. Gradients of the objective can be computed by automatic differentiation. That approach is often referred to as **discretize-then-optimize**. At first glance, this approach may suffer from very high memory requirements. Indeed, to get an accurate solution of the ODE, a numerical integration method may require K to be very large. Since a naive implementation of reverse-mode automatic differentiation has a memory that scales linearly with K , computing the gradient by a discretize-then-optimize method could be prohibitive. However, the memory requirements may easily be amortized using checkpointing, as explained in Section 8.5; see also (Gholaminejad *et al.*, 2019).

As for the **optimize-then-discretize** method, we still accumulate some truncation errors in the forward discretization process. This discretization error occurs when computing the gradient in reverse-mode too. The discretize-then-optimize method can be seen as computing gradients of a surrogate objective. For that objective, the gradients are correct and well-defined. However, they may not match the gradients of the true ODE formulation.

To compare the discretize-then-optimize and optimize-then-discretize approaches, Gholaminejad *et al.* (2019) compared their performance on an ODE whose solution can be computed analytically by selecting h to be linear in \mathbf{s} . The authors observed that discretize-then-optimize generally outperformed optimize-then-discretize. A middle ground can actually be found by using reversible differentiation schemes.

12.6.5 Reversible discretization schemes

Our exposition of the optimize-then-discretize or discretize-then-optimize approaches used a simple Euler explicit discretization scheme. However, for both approaches, we could have used other discretization schemes instead, such as reversible discretization schemes.

A reversible discretization scheme is a discretization scheme such that we have access to a closed-form formula for the inverse of its discretization step. Formally, a discretization method \mathcal{M} builds an approximation $(\mathbf{s}_k)_{k=1}^K$ of the solution of an ODE $\mathbf{s}'(t) = h(t, \mathbf{s}(t))$ on

an interval $[0, T]$ by computing for $k \in (1, \dots, K)$

$$t_k, \mathbf{s}_k, \mathbf{c}_k = \mathcal{M}(t_{k-1}, \mathbf{s}_{k-1}, \mathbf{c}_{k-1}; h, \delta), \quad (12.8)$$

where $\delta > 0$ is some fixed discretization step, t_k is the time step (typically $t_k = t_{k-1} + \delta$), \mathbf{s}_k is the approximation of $\mathbf{s}(t_k)$, and \mathbf{c}_k is some additional context variables used by the discretization method to build the iterates. An explicit Euler method does not have a context, but just as an optimization method may update some internal states, a discretization method can update some context variable. The discretization scheme in Eq. (12.8) is a forward discretization scheme as we took a positive discretization step. By taking a negative discretization step, we obtain the corresponding backward discretization scheme, for $k \in (K, \dots, 1)$,

$$t_{k-1}, \mathbf{s}_{k-1}, \mathbf{c}_{k-1} = \mathcal{M}(t_k, \mathbf{s}_k, \mathbf{c}_k; h, -\delta).$$

A discretization method is **reversible** if we have access to \mathcal{M}^{-1} to recompute the inputs of the discretization step from its outputs,

$$t_{k-1}, \mathbf{s}_{k-1}, \mathbf{c}_{k-1} = \mathcal{M}^{-1}(t_k, \mathbf{s}_k, \mathbf{c}_k; h, \delta).$$

A reversible discretization method is **symmetric** if the backward discretization scheme is exactly the inverse of the forward discretization scheme, i.e.,

$$\mathcal{M}(t_k, \mathbf{s}_k, \mathbf{c}_k; h, -\delta) = \mathcal{M}^{-1}(t_k, \mathbf{s}_k, \mathbf{c}_k; h, \delta).$$

The explicit Euler method is clearly not symmetric and a priori not reversible, unless we can solve for \mathbf{y}_{k-1} , the equation $\mathbf{y}_k = \mathbf{y}_{k-1} - \delta f(\mathbf{y}_{k-1})$.

Leapfrog method

The (asynchronous) **leapfrog method** (Zhuang *et al.*, 2021; Mutze, 2013) on the other hand is an example of symmetric reversible discretization method. For a constant discretization step δ , given $t_{k-1}, \mathbf{s}_{k-1}, \mathbf{c}_{k-1}$

and a function h , it computes

$$\begin{aligned}
 \bar{t}_{k-1} &:= t_{k-1} + \frac{\delta}{2} \\
 \bar{\mathbf{s}}_{k-1} &:= \mathbf{s}_{k-1} + \frac{\delta}{2} \mathbf{c}_{k-1} \\
 \bar{\mathbf{c}}_{k-1} &:= h(\bar{t}_{k-1}, \bar{\mathbf{s}}_{k-1}) \\
 t_k &:= \bar{t}_{k-1} + \frac{\delta}{2} \\
 \mathbf{s}_k &:= \bar{\mathbf{s}}_{k-1} + \frac{\delta}{2} \bar{\mathbf{c}}_{k-1} \\
 \mathbf{c}_k &:= 2\bar{\mathbf{c}}_{k-1} - \mathbf{c}_{k-1} \\
 \mathcal{M}(t_{k-1}, \mathbf{s}_{k-1}, \mathbf{c}_{k-1}; h, \delta) &:= (t_k, \mathbf{s}_k, \mathbf{c}_k).
 \end{aligned}$$

One can verify that we indeed have $\mathcal{M}(t_k, \mathbf{s}_k, \mathbf{c}_k; h, -\delta) = (t_{k-1}, \mathbf{s}_{k-1}, \mathbf{c}_{k-1})$.

By using a reversible symmetric discretization scheme in the optimize-then-discretize approach, we ensure that, at the end of the backward discretization pass, we recover exactly the original input. Therefore, by repeating forward and backward discretization schemes we always get the same gradient, which was not the case for an Euler explicit scheme.

By using a reversible discretization scheme in the discretize-then-optimize method, we address the memory issues of reverse mode autodiff. As explained in Section 8.6, we can recompute intermediate values during the backward pass rather than storing them.

Momentum residual networks

In the leapfrog method, the additional variables \mathbf{c}_k may actually be interpreted as velocities of a system whose acceleration is driven by the given function, that is, $\mathbf{s}''(t) = h(t, \mathbf{s}(t), \mathbf{w})$. Such an interpretation suggests alternatives to the usual neural ODE paradigm. For instance, **momentum neural networks** (Sander *et al.*, 2021b), can be interpreted as the discretization of a **second-order ordinary differential equations**, which are naturally amenable to reversible differentiation schemes with a low memory footprint.

12.6.6 Proof of the continuous adjoint method

In the following, we denote $\mathbf{s}(t, \mathbf{x}, \mathbf{w})$ the solution of the ODE at time t given the input \mathbf{x} and the parameters \mathbf{w} . We focus here on the formulation of the VJP. The proof relies on the existence of partial derivatives of $\mathbf{s}(t, \mathbf{x}, \mathbf{w})$, which we do not cover here and refer to, e.g., Pontryagin (1985) for a complete proof of such facts given the assumptions.

We use the ODE constraint to introduce adjoint variables, this time in the form of a continuously differentiable function \mathbf{r} . For any such a function \mathbf{r} , we have

$$\begin{aligned} \langle f(\mathbf{x}, \mathbf{w}), \mathbf{u} \rangle &= \langle \mathbf{s}(T, \mathbf{x}, \mathbf{w}), \mathbf{u} \rangle \\ &+ \int_0^T \langle \mathbf{r}(t), h(t, \mathbf{s}(t, \mathbf{x}, \mathbf{w}), \mathbf{w}) - \partial_t \mathbf{s}(t, \mathbf{x}, \mathbf{w}) \rangle dt, \end{aligned}$$

using Leibniz notations such as $\partial_t \mathbf{s}(t, \mathbf{x}, \mathbf{w}) = \partial_1 \mathbf{s}(t, \mathbf{x}, \mathbf{w})$. The VJPs then decompose as

$$\begin{aligned} &\partial_{\mathbf{w}} f(\mathbf{x}, \mathbf{w})^* [\mathbf{u}] \\ &= \partial_{\mathbf{w}} \mathbf{s}(T, \mathbf{x}, \mathbf{w})^* \mathbf{u} \\ &+ \int_0^T (\partial_{\mathbf{w}} \mathbf{s}(t, \mathbf{x}, \mathbf{w})^* \partial_{\mathbf{s}}^* h(t, \mathbf{s}(t, \mathbf{x}, \mathbf{w}), \mathbf{w})^* - \partial_{\mathbf{w}t}^2 \mathbf{s}(t, \mathbf{x}, \mathbf{w})^*) \mathbf{r}(t) dt \\ &+ \int_0^T \partial_{\mathbf{w}} h(t, \mathbf{s}(t, \mathbf{x}, \mathbf{w}), \mathbf{w})^* \mathbf{r}(t) dt, \\ &\partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{w})^* [\mathbf{u}] \\ &= \partial_{\mathbf{x}} \mathbf{s}(T, \mathbf{x}, \mathbf{w})^* \mathbf{u} \\ &+ \int_0^T (\partial_{\mathbf{x}} \mathbf{s}(t, \mathbf{x}, \mathbf{w})^* \partial_{\mathbf{s}}^* h(t, \mathbf{s}(t, \mathbf{x}, \mathbf{w}), \mathbf{w})^* - \partial_{\mathbf{x}t}^2 \mathbf{s}(t, \mathbf{x}, \mathbf{w})^*) \mathbf{r}(t) dt \end{aligned}$$

Here the second derivative terms $\partial_{\mathbf{w}t}^2 \mathbf{s}(t, \mathbf{x}, \mathbf{w})^* \mathbf{r}$, $\partial_{\mathbf{x}t}^2 \mathbf{s}(t, \mathbf{x}, \mathbf{w})^* \mathbf{r}$ correspond to second derivatives of $\langle \mathbf{s}(t, \mathbf{x}, \mathbf{w}), \mathbf{r} \rangle$. Since the Hessian is symmetric (Schwartz's theorem presented in Proposition 2.10), we can swap the derivatives in t and \mathbf{w} or \mathbf{x} . Then, to express the gradient uniquely in terms of first derivatives of \mathbf{s} , we use an integration by part

to have for example

$$\begin{aligned} \int_0^T \partial_{wt}^2 s(t, \mathbf{x}, \mathbf{w})^* \mathbf{r}(t) dt &= \int_0^T \partial_{tw}^2 s(t, \mathbf{x}, \mathbf{w})^* \mathbf{r}(t) dt \\ &= (\partial_w s(T, \mathbf{x}, \mathbf{w})^* \mathbf{r}(T) - \partial_w s(0, \mathbf{x}, \mathbf{w})^* \mathbf{r}(0)) \\ &\quad - \int_0^T \partial_w s(t, \mathbf{x}, \mathbf{w})^* \mathbf{r}(t)^* \partial_t \mathbf{r}(t) dt. \end{aligned}$$

Since $s(0) = \mathbf{x}$, we have $\partial_w s(0, \mathbf{x}, \mathbf{w})^* \mathbf{r}(0) = 0$. The VJP w.r.t. \mathbf{w} can then be written as

$$\begin{aligned} \partial_w f(\mathbf{x}, \mathbf{w})^* [\mathbf{u}] &= \partial_w s(T, \mathbf{x}, \mathbf{w})^* [\mathbf{u} - \mathbf{r}(T)] \\ &\quad + \int_0^T \partial_w s(t, \mathbf{x}, \mathbf{w})^* [\partial_s h(\mathbf{t}, s(t, \mathbf{x}, \mathbf{w}), \mathbf{w})^* \mathbf{r}(t) + \partial_t \mathbf{r}(t)] dt \\ &\quad + \int_0^T \partial_w h(t, s(t, \mathbf{x}, \mathbf{w}), \mathbf{w})^* \mathbf{r}(t) dt. \end{aligned}$$

By choosing $\mathbf{r}(t)$ to satisfy the adjoint ODE

$$\partial_t \mathbf{r}(t) = -\partial_s h(\mathbf{t}, s(t, \mathbf{x}, \mathbf{w}), \mathbf{w})^* \mathbf{r}(t), \quad \mathbf{r}(T) = \mathbf{u},$$

the expression of the VJP simplifies as

$$\partial_w f(\mathbf{x}, \mathbf{w})^* [\mathbf{u}] = \int_0^T \partial_w h(t, s(t, \mathbf{x}, \mathbf{w}), \mathbf{w})^* \mathbf{r}(t) dt.$$

For the VJP w.r.t. to \mathbf{x} , we can proceed similarly. Using an integration by part, we have, this time, $\partial_x s(0, \mathbf{x}, \mathbf{w})^* \mathbf{r}(0) = \mathbf{r}(0)$ since $s(0) = \mathbf{x}$. Choosing the same curve $\mathbf{r}(t)$ satisfying the adjoint ODE we get

$$\partial_x f(\mathbf{x}, \mathbf{w})^* [\mathbf{u}] = \mathbf{r}(0).$$

The existence of a curve \mathbf{r} solution of the backward ODE can easily be shown from Picard Lindelöf's theorem and the assumptions.

12.7 Summary

- We studied how to differentiate integrals, with a focus on expectations and solutions of a differential equation.

- For differentiating through expectations, we studied two main methods: the score function estimator (SFE, a.k.a. REINFORCE) and the path gradient estimator (PGE, a.k.a. reparametrization trick).
- The SFE is suitable when it is easy to sample from the distribution and its log-PDF is **explicitly** available. It is an unbiased estimator, but is known to suffer from high variance.
- The PGE is suitable for pushforward distributions, distributions that are **implicitly** defined through a transformation, or a sequence of them. These distributions can be easily sampled from, by injecting a source of randomness (such as noise) through the transformations. An unbiased, low-variance estimator of the gradient of their expectation is easily obtained, provided that we can interchange integration and differentiation.
- If we have an explicit distribution, we can sometimes convert it to an implicit distribution, thanks to the **location-scale transformation** or the **inverse transformation**.
- Conversely, if we have an implicit distribution, we can convert it to an explicit distribution using the **change-of-variables theorem**. However, this formula requires to compute the determinant of an inverse Jacobian, and is computationally expensive in general. Normalizing flows use invertible transformations so that the inverse Jacobian is cheap to compute, by design.
- **Stochastic computation graphs** can use a mix of explicit and implicit distributions at each node.
- For differentiating through the solution of a differential equation, two approaches can be considered.
- We can express the gradient as the solution of a differential equation thanks to the **continuous adjoint method**. We may then discretize backwards in time the differential equation that the gradient satisfies. This is the **optimize-then-discretize** approach.

- We can also first discretize the problem in such a way that the gradient can simply be computed by reverse mode auto-diff, applied on the discretization steps. This is the **discretize-then-optimize** approach. The optimize-then-discretize approach has no memory cost, but discrepancies between the forward and backward discretization passes often lead to numerical errors. The discretize-then-optimize introduces no such discrepancies but may come at a large memory cost.
- **Reversible discretization schemes** can circumvent the memory cost, as they enable the recomputation of intermediate discretization steps backwards in time.

Part IV

Smoothing programs

13

Smoothing by optimization

When a function is non-differentiable (or worse, discontinuous), a reasonable approach is to replace it by a differentiable approximation (or at least, by a continuous relaxation). We refer to the process of transforming a non-differentiable function into a differentiable one as “smoothing” the original function. In this chapter, we begin by reviewing a smoothing technique based on **infimal convolution**. We then review an equivalent dual approach, based on the **Legendre-Fenchel transform**. We illustrate how to apply these techniques to compute smoothed ReLUs and smoothed max operators, as well as continuous relaxations of step functions and argmax operators.

13.1 Primal approach

We first review how to smooth functions in the original, primal space of the function, using the infimal convolution and more particularly the Moreau envelope, a.k.a. Moreau-Yoshida regularization. In this chapter, we consider functions taking potentially infinite positive values, that is, functions taking values in the half-extended real line $\mathbb{R} \cup \{\infty\}$. For a

function $f: \mathbb{R}^M \rightarrow \mathbb{R} \cup \infty$, we define its domain as

$$\text{dom}(f) = \{\mathbf{u} \in \mathbb{R}^M : f(\mathbf{u}) < \infty\}.$$

13.1.1 Infimal convolution

Sometimes abbreviated inf-conv, the infimal convolution between two functions f and g creates a new function $f \square g$. It is defined as follows.

Definition 13.1 (Infimal convolution). The infimal convolution between two functions $f: \mathbb{R}^M \rightarrow \mathbb{R} \cup \{\infty\}$ and $g: \mathbb{R}^M \rightarrow \mathbb{R} \cup \{\infty\}$ is defined by

$$\begin{aligned} (f \square g)(\boldsymbol{\mu}) &:= \inf_{\mathbf{u} \in \mathbb{R}^M} f(\mathbf{u}) + g(\boldsymbol{\mu} - \mathbf{u}) \\ &= \inf_{\mathbf{z} \in \mathbb{R}^M} f(\boldsymbol{\mu} + \mathbf{z}) + g(\mathbf{z}) \\ &= \inf_{\mathbf{u}, \mathbf{z} \in \mathbb{R}^M} f(\mathbf{u}) + g(\mathbf{z}) \text{ s.t. } \mathbf{u} = \boldsymbol{\mu} + \mathbf{z}. \end{aligned}$$

It is easy to check that the three definitions are indeed equivalent, by using the change of variable $\mathbf{u} := \boldsymbol{\mu} + \mathbf{z}$, which is a location-scale transform; see Section 12.4.1.

The infimal convolution can be seen as a counterpart of the usual convolution, in which integration has been replaced by minimization (hence its name). Similarly to the classical convolution, it is **commutative**, meaning that for all $\boldsymbol{\mu} \in \mathbb{R}^M$, we have

$$(f \square g)(\boldsymbol{\mu}) = (g \square f)(\boldsymbol{\mu}).$$

Computing the infimal convolution involves the resolution of a minimization problem, that may or may not enjoy an analytical solution. Some examples are given in Table 13.1.

Existence

The infimal convolution $(f \square g)(\boldsymbol{\mu})$ exists if the infimum $\inf_{\mathbf{u} \in \mathbb{R}^M} f(\mathbf{u}) + g(\boldsymbol{\mu} - \mathbf{u})$ is finite (Bauschke and Combettes, 2017, Proposition 12.6). A sufficient condition to achieve this is that $\mathbf{u} \mapsto f(\mathbf{u}) + g(\boldsymbol{\mu} - \mathbf{u})$ is convex for all $\boldsymbol{\mu} \in \mathbb{R}^M$. However, this is not a necessary condition. For

Table 13.1: Examples of infimal convolutions. We use $\iota_{\mathcal{C}}$ to denote the indicator function of the set \mathcal{C} .

$f(\mathbf{u})$	$g(\mathbf{z})$	$(f \square g)(\boldsymbol{\mu})$
$f(\mathbf{u})$	0	$\inf_{\mathbf{u} \in \mathbb{R}^M} f(\mathbf{u})$
$f(\mathbf{u})$	$\iota_{\{\mathbf{v}\}}(\mathbf{z})$	$f(\boldsymbol{\mu} - \mathbf{v})$
$\iota_{\mathcal{C}}(\mathbf{u})$	$\iota_{\mathcal{D}}(\mathbf{z})$	$\iota_{\mathcal{C} + \mathcal{D}}(\boldsymbol{\mu})$
$\iota_{\mathcal{C}}(\mathbf{u})$	$\ \mathbf{z}\ _2$	$d_{\mathcal{C}}(\boldsymbol{\mu}) = \inf_{\mathbf{u} \in \mathcal{C}} \ \boldsymbol{\mu} - \mathbf{u}\ _2$
$f(\mathbf{u})$	$\frac{1}{2}\ \mathbf{z}\ _2^2$	$\text{env}_f(\boldsymbol{\mu}) = \inf_{\mathbf{u} \in \mathbb{R}^M} \frac{1}{2}\ \boldsymbol{\mu} - \mathbf{u}\ _2^2 + f(\mathbf{u})$

example, the infimum can be finite even if f or g are nonconvex, for example if their domain is a compact set.

Infimal convolution with a regularization function

When a function f is non-differentiable, a commonly-used technique is to replace it by its infimal convolution $f \square R$, with some regularization R . The most used regularization is the squared 2-norm, leading to the Moreau envelope, as we now review.

13.1.2 Moreau envelope

When $R(\mathbf{z}) := \frac{1}{2}\|\mathbf{z}\|_2^2$, the infimal convolution $f \square R$ gives the so-called **Moreau envelope** of f , which is also known as Moreau-Yoshida regularization of f .

Definition 13.2 (Moreau envelope). Given a function $f: \mathbb{R}^M \rightarrow \mathbb{R} \cup \{\infty\}$, its Moreau envelope is defined as

$$\begin{aligned}
 \text{env}_f(\boldsymbol{\mu}) &:= \left(f \square \frac{1}{2}\|\cdot\|_2^2 \right)(\boldsymbol{\mu}) \\
 &= \inf_{\mathbf{u} \in \mathbb{R}^M} f(\mathbf{u}) + \frac{1}{2}\|\boldsymbol{\mu} - \mathbf{u}\|_2^2 \\
 &= \inf_{\mathbf{z} \in \mathbb{R}^M} f(\boldsymbol{\mu} + \mathbf{z}) + \frac{1}{2}\|\mathbf{z}\|_2^2.
 \end{aligned}$$

Intuitively, the Moreau envelope is the minimal value over $\mathbf{u} \in \mathbb{R}^M$ of a trade-off between staying close to the input $\boldsymbol{\mu}$ according to the **proximity term** $\frac{1}{2}\|\boldsymbol{\mu} - \mathbf{u}\|_2^2$ and minimizing $f(\mathbf{u})$. Provided that the minimizer exists and is unique, we can define the associated **proximal operator** of f as

$$\text{prox}_f(\boldsymbol{\mu}) := \arg \min_{\mathbf{u} \in \mathbb{R}^M} \frac{1}{2}\|\boldsymbol{\mu} - \mathbf{u}\|_2^2 + f(\mathbf{u}),$$

In other words, we have for $\text{prox}_f(\boldsymbol{\mu})$ well defined,

$$\text{env}_f(\boldsymbol{\mu}) = f(\text{prox}_f(\boldsymbol{\mu})) + \frac{1}{2}\|\boldsymbol{\mu} - \text{prox}_f(\boldsymbol{\mu})\|_2^2. \quad (13.1)$$

Properties

A crucial property of the Moreau envelope env_f is that for any convex function f , it is always a smooth function, even when f itself is not smooth. By smooth, we formally mean that the resulting function env_f is differentiable everywhere with Lipschitz-continuous gradients. We say L -smooth, if the gradients are L -Lipschitz continuous. Such a property can determine the efficiency of optimization algorithms as reviewed in Section 15.4. We recap below useful properties of the Moreau envelope.

Proposition 13.1 (Properties of Moreau envelope). Let $f: \mathbb{R}^M \rightarrow \mathbb{R} \cup \{\infty\}$.

1. **Smoothness:** If f is convex, the function env_f is 1-smooth.
2. **Gradient:** Provided that $\text{prox}_f(\boldsymbol{\mu})$ is well-defined on $\boldsymbol{\mu} \in \mathbb{R}^M$, the gradient of the Moreau envelope can be expressed in terms of the proximal operator as

$$\nabla \text{env}_f(\boldsymbol{\mu}) = \boldsymbol{\mu} - \text{prox}_f(\boldsymbol{\mu}).$$

3. **Moreau decomposition:** If f is convex, then for any $\boldsymbol{\mu} \in \mathbb{R}^M$, we have the following identity

$$\text{prox}_f(\boldsymbol{\mu}) + \text{prox}_{f^*}(\boldsymbol{\mu}) = \boldsymbol{\mu},$$

where f^* is the convex conjugate of f , detailed in Section 13.2. In particular, we get

$$\nabla \text{env}_f(\boldsymbol{\mu}) = \text{prox}_f(\boldsymbol{\mu})$$

4. **Convexity:** env_f is convex if f is convex.
5. **Infimums coincide** env_f has the same infimum as the original function f :

$$\min_{\boldsymbol{\mu} \in \mathbb{R}^M} \text{env}_f(\boldsymbol{\mu}) = \min_{\mathbf{u} \in \mathbb{R}^M} f(\mathbf{u}).$$

Proof.

1. This is best seen using the dual approach detailed in Section 13.3.
2. This follows from Danskin's theorem, reviewed in Section 11.2.
3. See, e.g., Bauschke and Combettes (2017, Theorem 14.3).
4. This follows from the fact that the infimum of a jointly convex function is convex.
5. We have

$$\begin{aligned} \inf_{\boldsymbol{\mu} \in \mathbb{R}^M} \text{env}_f(\boldsymbol{\mu}) &= \inf_{\boldsymbol{\mu} \in \mathbb{R}^M} \inf_{\mathbf{u} \in \mathbb{R}^M} \frac{1}{2} \|\boldsymbol{\mu} - \mathbf{u}\|_2^2 + f(\mathbf{u}) \\ &= \inf_{\mathbf{u} \in \mathbb{R}^M} \inf_{\boldsymbol{\mu} \in \mathbb{R}^M} \frac{1}{2} \|\boldsymbol{\mu} - \mathbf{u}\|_2^2 + f(\mathbf{u}) \\ &= \inf_{\mathbf{u} \in \mathbb{R}^M} f(\mathbf{u}). \end{aligned}$$

□

Examples

To illustrate smoothing from the Moreau envelope perspective, we show how to smooth the 1-norm. In this case, we obtain an analytical expression for the Moreau envelope.

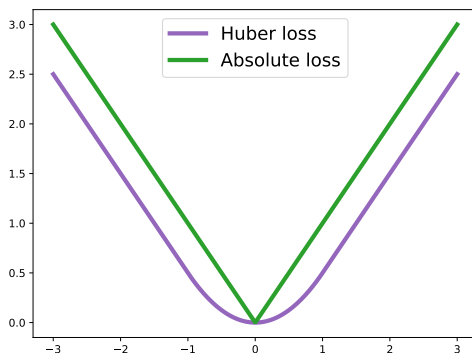


Figure 13.1: The Huber loss is the Moreau envelope of the absolute loss.

Example 13.1 (Smoothing the 1-norm via infimal convolution). We wish to smooth $f(\mathbf{u}) := \|\mathbf{u}\|_1 = \sum_{j=1}^M |u_j|$. The corresponding proximal operator is the soft-thresholding operator (see Section 16.4),

$$\begin{aligned} \text{prox}_f(\boldsymbol{\mu}) &= \arg \min_{\mathbf{u} \in \mathbb{R}^M} \frac{1}{2} \|\boldsymbol{\mu} - \mathbf{u}\|_2^2 + \|\mathbf{u}\|_1 \\ &= \text{sign}(\boldsymbol{\mu}) \cdot \max(|\boldsymbol{\mu}| - 1, 0). \end{aligned}$$

Using Eq. (13.1) and after some algebraic manipulations, we obtain

$$\text{env}_f(\boldsymbol{\mu}) = \sum_{j=1}^M \text{huber}(\mu_j) \approx \sum_{j=1}^M |\mu_j|,$$

where we defined the **Huber loss**

$$\text{huber}(\mu_j) := \begin{cases} \frac{\mu_j^2}{2} & \text{if } |\mu_j| \leq 1 \\ |\mu_j| - \frac{1}{2} & \text{if } |\mu_j| > 1 \end{cases}.$$

This is illustrated in Fig. 13.1 with $M = 1$.

We also illustrate in Fig. 13.2 that the Moreau envelope of nonconvex functions can be approximately computed numerically.

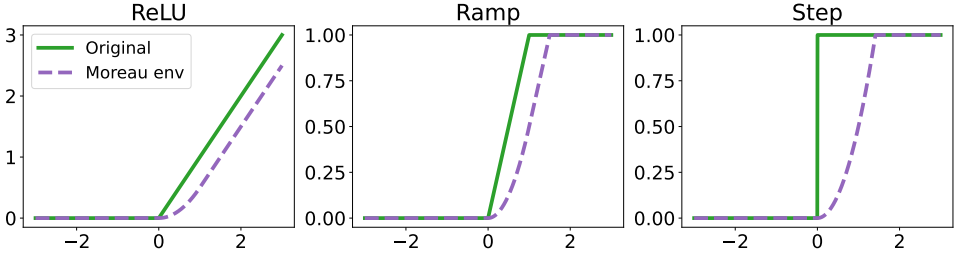


Figure 13.2: The Moreau envelope is not limited to convex functions. For instance, the ramp function is continuous but nonconvex, and the step function is not only nonconvex but also discontinuous. In this figure, we approximately computed the infimum over $u \in \mathbb{R}$ in Definition 13.2 by restricting the search on a finite grid, in a closed interval.

13.1.3 Vector-valued functions

The Moreau envelope is defined by $\text{env}_f(\boldsymbol{\mu}) := \inf_{\mathbf{u} \in \mathbb{R}^M} f(\mathbf{u}) + \frac{1}{2} \|\boldsymbol{\mu} - \mathbf{u}\|_2^2$. As such, it is limited to scalar-valued functions $f: \mathbb{R}^M \rightarrow \mathbb{R}$. To extend the Moreau envelope to vector-valued functions $f: \mathbb{R}^M \rightarrow \mathbb{R}^T$, where $f(\mathbf{u}) = (f_1(\mathbf{u}), \dots, f_T(\mathbf{u}))$ and $f_i: \mathbb{R}^M \rightarrow \mathbb{R}$ for $i \in [T]$, we may choose to smooth each f_j separately to define

$$\mathbf{env}_f(\boldsymbol{\mu}) := (\text{env}_{f_1}(\boldsymbol{\mu}), \dots, \text{env}_{f_T}(\boldsymbol{\mu})),$$

where

$$\text{env}_{f_i}(\boldsymbol{\mu}) = \inf_{\mathbf{u}_i \in \mathbb{R}^M} f_i(\mathbf{u}_i) + \frac{1}{2} \|\boldsymbol{\mu} - \mathbf{u}_i\|_2^2.$$

This approach requires to solve T separate minimization problems and performs the smoothing of each output coordinate $i \in [T]$ independently. From Proposition 2.9, we then have that the VJP of $\mathbf{env}_f(\mathbf{u})$ with any direction $\mathbf{d} \in \mathbb{R}^T$ is

$$\begin{aligned} \partial \mathbf{env}_f(\boldsymbol{\mu})^*[\mathbf{d}] &= \sum_{i=1}^T \partial \text{env}_{f_i}(\boldsymbol{\mu})^*[d_i] \\ &= \sum_{i=1}^T d_i \nabla \text{env}_{f_i}(\boldsymbol{\mu}). \end{aligned}$$

In the particular case $f(\mathbf{u}) = (f_1(u_1), \dots, f_T(u_T))$, we obtain

$$\partial \mathbf{env}_f(\boldsymbol{\mu})^*[\mathbf{d}] = \sum_{i=1}^T d_i \text{env}_{f_i}(\mu_i).$$

An alternative was proposed by Roulet and Harchaoui (2022). For a differentiable function $f: \mathbb{R}^M \rightarrow \mathbb{R}^T$, we recall that the VJP of f with a direction $\mathbf{d} \in \mathbb{R}^T$ reads

$$\partial f(\mathbf{u})^*[\mathbf{d}] = \nabla \langle f, \mathbf{d} \rangle(\mathbf{u}),$$

where we defined the scalar-valued function $\langle f, \mathbf{d} \rangle(\mathbf{u}) := \langle f(\mathbf{u}), \mathbf{d} \rangle$. As a result, if f is non differentiable, a natural idea is to approximate its VJP $\partial f(\mathbf{u})^*[\mathbf{d}]$ (had it existed) by the gradient $\nabla \text{env}_{\langle f, \mathbf{d} \rangle}(\boldsymbol{\mu})$ of the Moreau envelope

$$\text{env}_{\langle f, \mathbf{d} \rangle}(\boldsymbol{\mu}) = \inf_{\mathbf{u} \in \mathbb{R}^M} \langle f(\mathbf{u}), \mathbf{d} \rangle + \frac{1}{2} \|\boldsymbol{\mu} - \mathbf{u}\|_2^2. \quad (13.2)$$

This requires a single optimization problem to solve, independently of the number of outputs T . Moreover, for $\mathbf{d} = \mathbf{e}_i$, this recovers $\text{env}_{f_i}(\boldsymbol{\mu})$ as a special case.

This approach allows in principle to perform reverse-mode autodiff (gradient backpropagation) on a neural network whose layers use the Moreau envelope. Indeed, following Proposition 13.1, the approximate VJP of f with a direction \mathbf{d} is given by

$$\partial f(\boldsymbol{\mu})^*[\mathbf{d}] \approx \nabla \text{env}_{\langle f, \mathbf{d} \rangle}(\boldsymbol{\mu}) = \boldsymbol{\mu} - \mathbf{u}^*,$$

where \mathbf{u}^* is the solution of the minimization problem in Eq. (13.2). However, we emphasize that this minimization problem could be difficult to solve in general. Indeed, when performing gradient backpropagation, the direction \mathbf{d} is not necessarily non-negative, therefore the function being minimized in Eq. (13.2) could be nonconvex, even if each f_i is convex. Another potential caveat is that the direction \mathbf{d} influences the smoothing strength, while in principle we should be able to smooth a function independently of whether we compute its VJP or not. To see that, for example in the particular case $f(\mathbf{u}) = (f_1(u_1), \dots, f_T(u_T))$, one easily checks that for $\mathbf{d} = (d_1, \dots, d_T)$, we get

$$\text{env}_{\langle f, \mathbf{d} \rangle}(\boldsymbol{\mu}) = \sum_{i=1}^T \text{env}_{d_i f_i}(\mu_i).$$

Smoothing vector-valued functions by Moreau envelope (or more generally, by infimal convolution) remains an open area of research. We will see in Chapter 14 that smoothing by convolution more naturally supports vector-valued functions.

13.2 Legendre–Fenchel transforms, convex conjugates

The Legendre–Fenchel transform, a.k.a. convex conjugate, is a way to turn a function f into a new function, denoted f^* . We now review it in detail, as it plays a major role for the dual approach to smoothing.

13.2.1 Definition

Consider the class of affine functions of the form

$$\mathbf{u} \mapsto \langle \mathbf{u}, \mathbf{v} \rangle - b.$$

These functions are parametrized by their slope $\mathbf{v} \in \mathbb{R}^M$ and their intercept $-b \in \mathbb{R}$. Now, suppose we fix \mathbf{v} . Given a function $f(\mathbf{u})$, affine lower bounds of $f(\mathbf{u})$ are all the functions of \mathbf{u} such that b satisfies for all $\mathbf{u} \in \mathbb{R}^M$,

$$\langle \mathbf{u}, \mathbf{v} \rangle - b \leq f(\mathbf{u}) \iff \langle \mathbf{u}, \mathbf{v} \rangle - f(\mathbf{u}) \leq b.$$

The **tightest** lower bound is then defined by b such that

$$b := \sup_{\mathbf{u} \in \text{dom}(f)} \langle \mathbf{u}, \mathbf{v} \rangle - f(\mathbf{u}),$$

where we recall that the domain of f is defined by

$$\text{dom}(f) := \{\mathbf{u} \in \mathbb{R}^M : f(\mathbf{u}) < \infty\}.$$

This leads to the definition of **Legendre–Fenchel transform**, a.k.a. **convex conjugate**.

Definition 13.3 (Legendre–Fenchel transform, convex conjugate). Given a function $f: \mathbb{R}^M \rightarrow \mathbb{R} \cup \{\infty\}$, its convex conjugate is defined by

$$f^*(\mathbf{v}) := \sup_{\mathbf{u} \in \text{dom}(f)} \langle \mathbf{u}, \mathbf{v} \rangle - f(\mathbf{u}).$$

We use a sup rather than a max to indicate that $f^*(\mathbf{v})$ is potentially ∞ . Following the previous discussion, $-f^*(\mathbf{v})$ is the intercept of the tightest affine lower bound with slope \mathbf{v} of $f(\mathbf{u})$. This is illustrated Fig. 13.3.

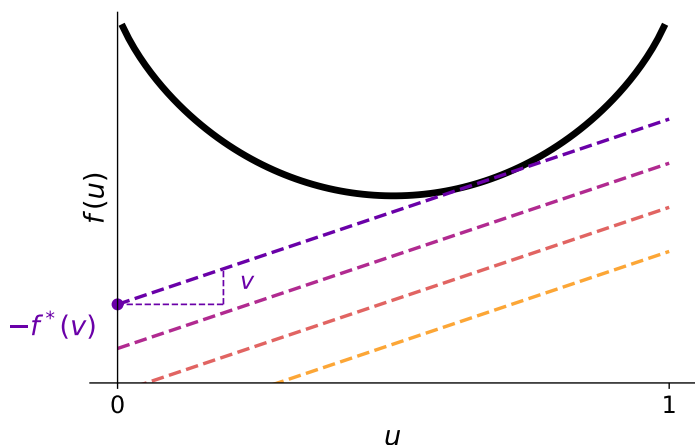


Figure 13.3: For a fixed slope v , the function $u \mapsto uv - f^*(v)$ is the highest affine lower bound of f with slope v .

The Legendre-Fenchel transform is a function transformation, as it produces a new function f^* . It can be seen as a dual representation of a function: instead of representing a convex function f by its graph $(u, f(u))$ for $u \in \text{dom}(f)$, we can represent it by the set of tangents with slope v and intercept $-f^*(v)$ for $v \in \text{dom}(f^*)$, as illustrated in Fig. 13.4. As the name “convex conjugate” indicates, it is convex, even if the original function is not.

13.2.2 Closed-form examples

Computing $f^*(v)$ involves the resolution of a maximization problem, which could be difficult in general without assumption on f . In some cases, however, we can compute an analytical expression, as we now illustrate.

Example 13.2 (Analytical conjugate examples). When $f(u) = \frac{1}{2}\|u\|_2^2$, with $\text{dom}(f) = \mathbb{R}^M$, the conjugate is

$$f^*(v) = \max_{u \in \mathbb{R}^M} \langle u, v \rangle - \frac{1}{2}\|u\|_2^2.$$

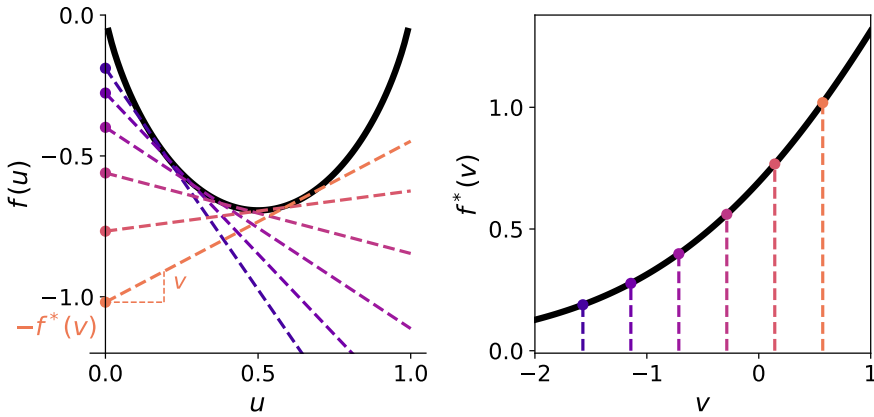


Figure 13.4: **Left:** instead of representing a convex function f by its graph $(u, f(u))$ for $u \in \text{dom}(f)$, we can represent it by the set of tangents with slope v and intercept $-f^*(v)$ for $v \in \text{dom}(f^*)$. **Right:** by varying the slope v of all possible tangents, we obtain a function of the slope v rather than of the original input u . The colors of the tangents on the left are chosen to match the colors of the vertical lines on the right.

Setting the gradient $u \mapsto \langle u, v \rangle - \frac{1}{2}\|u\|_2^2$ to zero, we obtain $u^* = v$. Plugging u^* back, we therefore obtain

$$f^*(v) = \langle u^*, v \rangle - \frac{1}{2}\|u^*\|_2^2 = \frac{1}{2}\|v\|_2^2.$$

Therefore, $f = f^*$ in this case.

When $f(u) = \langle u, \log u \rangle$, with $\text{dom}(f) = \mathbb{R}_+^M$, the minimizer of $u \mapsto \langle u, v \rangle - \langle u, \log u \rangle$ is $u^* = \exp(v - \mathbf{1})$ and the conjugate is

$$f^*(v) = \sum_{j=1}^M \exp(v_j - 1).$$

See for instance Boyd and Vandenberghe (2004) or Beck (2017) for many more examples.

Constraining the domain

We can incorporate constraints using an **indicator function** with values in the extended real line $\mathbb{R} \cup \{\infty\}$,

$$\iota_{\mathcal{C}}(\mathbf{u}) := \begin{cases} 0 & \text{if } \mathbf{u} \in \mathcal{C} \\ +\infty & \text{otherwise} \end{cases}.$$

Example 13.3 (Incorporating constraints). If $f(\mathbf{u}) = \iota_{\mathcal{C}}(\mathbf{u})$, where \mathcal{C} is a convex set, then

$$f^*(\mathbf{v}) = \sup_{\mathbf{u} \in \text{dom}(f)} \langle \mathbf{u}, \mathbf{v} \rangle - f(\mathbf{u}) = \sup_{\mathbf{u} \in \mathcal{C}} \langle \mathbf{u}, \mathbf{v} \rangle := \sigma_{\mathcal{C}}(\mathbf{v}),$$

which is known as the **support function** of \mathcal{C} . The corresponding argmax (assuming that it exists),

$$\mathbf{v} \mapsto \arg \max_{\mathbf{u} \in \mathcal{C}} \langle \mathbf{u}, \mathbf{v} \rangle,$$

is known as the **linear maximization oracle** (LMO) of \mathcal{C} . As another example, if $f(\mathbf{u}) = \langle \mathbf{u}, \log \mathbf{u} \rangle + \iota_{\Delta^M}(\mathbf{u})$ then

$$f^*(\mathbf{v}) = \text{logsumexp}(\mathbf{v}) = \log \sum_{i=1}^M \exp(v_j).$$

We postpone a proof to Proposition 13.9.

13.2.3 Properties

The conjugate enjoys several useful properties, that we now summarize.

Proposition 13.2 (Convex conjugate properties).

1. **Convexity:** $f^*(\mathbf{v})$ is a **convex** function for **all** $f: \mathbb{R}^M \rightarrow \mathbb{R} \cup \{\infty\}$ (even if f is nonconvex).
2. **Fenchel-Young inequality:** for all $\mathbf{u}, \mathbf{v} \in \mathbb{R}^M$

$$f(\mathbf{u}) + f^*(\mathbf{v}) - \langle \mathbf{u}, \mathbf{v} \rangle \geq 0.$$

3. **Gradient:** if the supremum in Definition 13.3 is uniquely achieved, then $f^*(\mathbf{v})$ is differentiable at \mathbf{v} and its gradient is

$$\nabla f^*(\mathbf{v}) = \arg \max_{\mathbf{u} \in \text{dom}(f)} \langle \mathbf{u}, \mathbf{v} \rangle - f(\mathbf{u}).$$

Otherwise, $f^*(\mathbf{v})$ is sub-differentiable at \mathbf{v} and we get a sub-gradient instead.

4. **Maps:** If f and f^* are differentiable, then

$$\mathbf{v} = \nabla f(\mathbf{u}) \iff \mathbf{u} = \nabla f^*(\mathbf{v}) \iff f^*(\mathbf{v}) + f(\mathbf{u}) - \langle \mathbf{u}, \mathbf{v} \rangle = 0.$$

5. **Biconjugate:** $f = f^{**}$ if and only if f is convex and closed (i.e., its sublevel sets form a closed set), otherwise $f^{**} \leq f$.

Proof.

1. This follows from the fact that $\mathbf{v} \mapsto \sup_{\mathbf{u} \in C} g(\mathbf{u}, \mathbf{v})$ is convex if g is convex in \mathbf{v} . Note that this is true even if g is nonconvex in \mathbf{u} . Here, $g(\mathbf{u}, \mathbf{v}) = \langle \mathbf{u}, \mathbf{v} \rangle - f(\mathbf{u})$, which is affine in \mathbf{v} and therefore convex in \mathbf{v} .
2. This follows immediately from Definition 13.3.
3. This follows from Danskin's theorem, reviewed in Section 11.2. Another way to see this is by observing that

$$\begin{aligned} f^*(\mathbf{v}) &= \langle \mathbf{g}, \mathbf{v} \rangle - f(\mathbf{g}) \\ f^*(\mathbf{v}') &\geq \langle \mathbf{g}, \mathbf{v}' \rangle - f(\mathbf{g}), \end{aligned}$$

where $\mathbf{g} := \arg \max_{\mathbf{u} \in \text{dom}(f)} \langle \mathbf{u}, \mathbf{v} \rangle - f(\mathbf{u})$. Subtracting the two, we obtain

$$f^*(\mathbf{v}') \geq f^*(\mathbf{v}) + \langle \mathbf{g}, \mathbf{v}' - \mathbf{v} \rangle.$$

Now, using that f^* is convex and Definition 15.6, we obtain that $\mathbf{g} = \nabla f^*(\mathbf{v})$.

4. See, e.g., Bauschke and Combettes (2017, Proposition 16.10).
5. See Boyd and Vandenberghe (2004, Section 3.3).

□

13.2.4 Conjugate calculus

While deriving a convex conjugate expression can be difficult in general, in some cases, it is possible to use simple rules to derive conjugates in terms of other conjugates.

Proposition 13.3 (Conjugate calculus rules).

1. **Separable sum of functions:** if $f(\mathbf{u}) = \sum_{j=1}^M f_j(u_j)$, then

$$f^*(\mathbf{v}) = \sum_{j=1}^M f_j^*(v_j).$$

2. **Scalar multiplication:** if $f(\mathbf{u}) = c \cdot g(\mathbf{u})$, for $c > 0$, then

$$f^*(\mathbf{v}) = c \cdot g^*(\mathbf{v}/c).$$

3. **Addition to an affine function and translation:** if $f(\mathbf{u}) = g(\mathbf{u}) + \langle \boldsymbol{\alpha}, \mathbf{u} \rangle + \beta$, then

$$f^*(\mathbf{v}) = g^*(\mathbf{v} - \boldsymbol{\alpha}) - \beta.$$

4. **Composition with an invertible linear map:** if $f(\mathbf{u}) = g(M\mathbf{u})$, where $x \mapsto Mx$ is an invertible linear map, then

$$f^*(\mathbf{v}) = g^*(M^{-T}\mathbf{v}).$$

5. **Non-separable sum of functions:** if h_1 and h_2 are convex functions, then $(h_1 + h_2)^* = h_1^* \square h_2^*$, where \square is the infimal convolution operator.

13.2.5 Fast Legendre transform

When an analytical expression is not available, we can resort to numerical schemes to approximately compute the transform / conjugate. When f is convex, because $-f$ is concave, the maximization in Definition 13.3 is that of a concave function. Therefore, the conjugate can be computed to arbitrary precision in polynomial time using classical iterative algo-

gorithms for constrained optimization such as projected gradient descent (Section 16.3) or conditional gradient a.k.a. Frank-Wolfe (Jaggi, 2013). Without convexity assumption on f , $f^*(\mathbf{v})$ can be approximated by

$$f^*(\mathbf{v}) \approx \sup_{\mathbf{u} \in \mathcal{U}} \langle \mathbf{u}, \mathbf{v} \rangle - f(\mathbf{u}),$$

where $\mathcal{U} \subseteq \text{dom}(f)$ is a discrete grid of values. We can then compute $f^*(\mathbf{v})$ for several inputs $\mathbf{v} \in \mathcal{V}$ using the linear-time Legendre transform algorithm (Lucet, 1997), where $\mathcal{V} \subseteq \text{dom}(f^*)$ is another discrete grid. The complexity is $O(|\mathcal{U}| \cdot |\mathcal{V}|)$, which is linear in the grid sizes. However, the grid sizes are typically $|\mathcal{U}| = |\mathcal{V}| = O(N^M)$, for N equally-distributed points in each of the M dimensions. Therefore, this approach is limited to small-dimensional settings, e.g., $M \in \{1, 2, 3\}$.

13.3 Dual approach

Previously, we presented how to smooth a function by performing its infimal convolution with a primal-space regularization R . We now present how to smooth a function by regularizing its Legendre-Fenchel transform (convex conjugate) instead. This dual, equivalent approach, is often mathematically more convenient.

13.3.1 Duality between strong convexity and smoothness

We begin by stating a well-known result that will underpin this whole section: smoothness and strong convexity are dual to each other (Hiriart-Urruty and Lemaréchal, 1993; Kakade *et al.*, 2009; Beck, 2017; Zhou, 2018).

Proposition 13.4 (Duality between strong convexity and smoothness).

f is $\frac{1}{\mu}$ -strongly convex w.r.t. the norm $\|\cdot\|$ over $\text{dom}(f)$ if and only if f^* is μ -smooth w.r.t. the dual norm $\|\cdot\|_*$ over $\text{dom}(f^*)$.

For a review of the notions of smoothness and strong convexity, see Section 15.4. We give two examples of strongly-convex and smooth conjugate pairs in Table 13.2.

Table 13.2: Examples of strongly-convex and smooth conjugate pairs.

Function	Norm	Domain	Conjugate	Dual norm	Dual domain
$\frac{1}{2}\ \mathbf{u}\ _2^2$	$\ \cdot\ _2$	\mathbb{R}^M	$\frac{1}{2}\ \mathbf{v}\ _2^2$	$\ \cdot\ _2$	\mathbb{R}^M
$\langle \mathbf{u}, \log \mathbf{u} \rangle$	$\ \cdot\ _1$	\triangle^M	$\text{logsumexp}(\mathbf{v})$	$\ \cdot\ _\infty$	\mathbb{R}^M

13.3.2 Smoothing by dual regularization

The duality between smoothness and strong convexity suggests a generic approach in order to smooth a function $f: \mathbb{R}^M \rightarrow \mathbb{R}$, by going through the **dual** space.

1. Compute the conjugate f^* :

$$f^*(\mathbf{v}) := \sup_{\mathbf{u} \in \text{dom}(f)} \langle \mathbf{u}, \mathbf{v} \rangle - f(\mathbf{u}).$$

2. Add strongly-convex regularization Ω to the conjugate:

$$f_\Omega^*(\mathbf{v}) := f^*(\mathbf{v}) + \Omega(\mathbf{v}). \quad (13.3)$$

3. Go back to the primal space, by computing the conjugate of f_Ω^* :

$$f_\Omega(\mathbf{u}) := f_\Omega^{**}(\mathbf{u}) = \max_{\mathbf{v} \in \mathbb{R}^M} \langle \mathbf{u}, \mathbf{v} \rangle - f_\Omega^*(\mathbf{v}).$$

Note that \mathbf{u} and \mathbf{v} belong to different spaces, i.e., $\mathbf{u} \in \text{dom}(f)$ and $\mathbf{v} \in \text{dom}(f^*)$. Following Proposition 13.4, if Ω is μ -strongly convex, then $f_\Omega(\mathbf{u})$ is $\frac{1}{\mu}$ -smooth. Furthermore, following Proposition 13.2, $f_\Omega(\mathbf{u})$ is convex, even if f is nonconvex. Therefore, $f_\Omega(\mathbf{u})$ is a **smooth and convex relaxation** of $f(\mathbf{u})$.

Steps 1 and 3 are the most challenging, as they both require the derivation of a conjugate. While an analytical solution may not exist in general, in some simple cases, there is, as we now illustrate.

Example 13.4 (Smoothing the 1-norm via dual regularization). We revisit Example 13.1, this time from the dual perspective. We wish to smooth out the 1-norm $f(\mathbf{u}) := \|\mathbf{u}\|_1 = \sum_{j=1}^M |u_j|$.

1. **Compute the conjugate.** The conjugate of any norm $\|\cdot\|$

is the indicator function of the dual norm's unit ball $\{\mathbf{v} \in \mathbb{R}^M : \|\mathbf{v}\|_* \leq 1\}$ (see e.g. Boyd and Vandenberghe (2004, Example 3.26)). The dual norm of $\|\mathbf{u}\|_1$ is $\|\mathbf{v}\|_\infty$. Moreover,

$$\{\mathbf{v} \in \mathbb{R}^M : \|\mathbf{v}\|_\infty \leq 1\} = [-1, 1]^M.$$

Recalling that $\iota_{\mathcal{C}}$ is the indicator function of \mathcal{C} , we obtain

$$f^*(\mathbf{v}) = \iota_{[-1,1]^M}(\mathbf{v}).$$

2. Adding strongly-convex regularization. We add quadratic regularization $\Omega(\mathbf{v}) := \frac{1}{2}\|\mathbf{v}\|_2^2$ to define

$$f_\Omega^*(\mathbf{v}) := \iota_{[-1,1]^M}(\mathbf{v}) + \Omega(\mathbf{v}).$$

3. Going back to the primal.

$$f_\Omega(\mathbf{u}) = f_\Omega^{**}(\mathbf{u}) = \langle \mathbf{u}, \mathbf{v}^* \rangle - \Omega(\mathbf{v}^*) = \sum_{i=1}^M \text{huber}(u_i),$$

where $\mathbf{v}^* = \text{clip}(\mathbf{u}) := \max(\min(\mathbf{u}, 1), -1)$.

We therefore indeed recover the Huber loss from Example 13.1.

ReLU functions can be smoothed out in a similar way, as we see in more details in Section 13.4.

The dual approach allows us to easily bound the smoothed function in terms of the original function.

Proposition 13.5 (Bounds). If $\mathcal{L}_\Omega \leq \Omega(\mathbf{v}) \leq \mathcal{U}_\Omega$ for all $\mathbf{v} \in \text{dom}(\Omega)$, then for all $\mathbf{u} \in \mathbb{R}^M$,

$$f(\mathbf{u}) - \mathcal{U}_\Omega \leq f_\Omega(\mathbf{u}) \leq f(\mathbf{u}) - \mathcal{L}_\Omega.$$

Proof. Let us define

$$\mathbf{v}^* := \arg \max_{\mathbf{v} \in \mathbb{R}^M} \langle \mathbf{u}, \mathbf{v} \rangle - f^*(\mathbf{v})$$

$$\mathbf{v}_\Omega^* := \arg \max_{\mathbf{v} \in \mathbb{R}^M} \langle \mathbf{u}, \mathbf{v} \rangle - f_\Omega^*(\mathbf{v}),$$

where we recall that $f_\Omega^* := f^* + \Omega$. We then have for all $\mathbf{u} \in \mathbb{R}^M$

$$f_\Omega(\mathbf{u}) = \langle \mathbf{u}, \mathbf{v}_\Omega^* \rangle - f_\Omega^*(\mathbf{v}_\Omega^*) \geq \langle \mathbf{u}, \mathbf{v}^* \rangle - f_\Omega^*(\mathbf{v}^*) = f(\mathbf{u}) - \Omega(\mathbf{v}^*)$$

and similarly

$$f(\mathbf{u}) - \Omega(\mathbf{v}_\Omega^*) = \langle \mathbf{u}, \mathbf{v}^* \rangle - f^*(\mathbf{v}^*) - \Omega(\mathbf{v}_\Omega^*) \geq \langle \mathbf{u}, \mathbf{v}_\Omega^* \rangle - f_\Omega^*(\mathbf{v}_\Omega^*) = f_\Omega(\mathbf{u}).$$

Combining the two with $\mathcal{L}_\Omega \leq \Omega(\mathbf{v}) \leq \mathcal{U}_\Omega$ for all $\mathbf{v} \in \text{dom}(\Omega)$, we obtain

$$f(\mathbf{u}) - \mathcal{U}_\Omega \leq f(\mathbf{u}) - \Omega(\mathbf{v}^*) \leq f_\Omega(\mathbf{u}) \leq f(\mathbf{u}) - \Omega(\mathbf{v}_\Omega^*) \leq f(\mathbf{u}) - \mathcal{L}_\Omega.$$

□

Remark 13.1 (The gradient is differentiable almost everywhere). From Proposition 13.2, the gradient of $f_\Omega(\mathbf{u})$ equals

$$\nabla f_\Omega(\mathbf{u}) = \arg \max_{\mathbf{v} \in \mathbb{R}^M} \langle \mathbf{u}, \mathbf{v} \rangle - f_\Omega^*(\mathbf{v}).$$

If Ω is strongly convex, then f_Ω is smooth, meaning that ∇f_Ω is Lipschitz continuous. From Rademacher's theorem reviewed in Section 2.7.1, ∇f_Ω is then differentiable almost everywhere (that is, f_Ω is twice differentiable almost everywhere). We use this property in the sequel to define continuous differentiable almost everywhere relaxations of step functions and argmax operators.

13.3.3 Equivalence between primal and dual regularizations

So far, we saw two approaches to obtain a smooth approximation of a function f . The first approach is based on the infimal convolution $f \square R$, where $R: \text{dom}(f) \rightarrow \mathbb{R}$ denotes primal regularization. The second approach is based on regularizing the Legendre-Fenchel transform (convex conjugate) f^* of f with some dual regularization Ω , to define $f_\Omega = (f^* + \Omega)^*$. It turns out that both approaches are equivalent.

Proposition 13.6 (Equivalence between primal and dual regularizations).

Let $f: \mathbb{R}^M \rightarrow \mathbb{R} \cup \{\infty\}$ and $R: \mathbb{R}^M \rightarrow \mathbb{R} \cup \{\infty\}$, both convex and closed. Then, $f_\Omega = (f^* + \Omega)^* = f \square R$ with $\Omega = R^*$.

Proof. We have

$$f_{\Omega}(\mathbf{u}) = (f^* + \Omega)^*(\mathbf{u}) = \sup_{\mathbf{v} \in \text{dom}(f^*)} \langle \mathbf{u}, \mathbf{v} \rangle - f^*(\mathbf{v}) - \Omega(\mathbf{v}).$$

If h_1 and h_2 are convex, we have $(h_1 + h_2)^* = h_1^* \square h_2^*$ (Beck, 2017, Theorem 4.17). Using $h_1 = f^*$ and $h_2 = \Omega = R^*$ gives the desired result using that $f^{**} = f$ and $R^{**} = R$ since both are convex and closed (see Proposition 13.2). \square

In particular, with $\Omega = \frac{1}{2} \|\cdot\|_2^2 = \Omega^*$, this shows that the Moreau envelope can equivalently be written as

$$\text{env}_f = f_{\Omega} = f_{\Omega^*}.$$

Given the equivalence between the primal and dual approaches, using one approach or the other is mainly a matter of mathematical or algorithmic convenience, depending on the case.

In this book, we focus on applications of smoothing techniques to differentiable programming. For applications to non-smooth optimization, see Nesterov (2005) and Beck and Teboulle (2012).

13.3.4 Regularization scaling

Dual approach

If Ω is 1-strongly convex, then f_{Ω} is a 1-smooth approximation of the original function f . To control the smoothness of the approximation, it suffices to regularize with $\varepsilon\Omega$ for $\varepsilon > 0$, leading to a $1/\varepsilon$ -smooth approximation $f_{\varepsilon\Omega}$ of f . Moreover, one can check that

$$\begin{aligned} f_{\varepsilon\Omega}(\mathbf{v}) &= \varepsilon f_{\Omega}(\mathbf{v}/\varepsilon) \\ \nabla f_{\varepsilon\Omega}(\mathbf{v}) &= \nabla f_{\Omega}(\mathbf{v}/\varepsilon). \end{aligned}$$

Therefore, if we know how to compute f_{Ω} , we can also compute $f_{\varepsilon\Omega}$ and its gradient easily. Furthermore, the approximation error induced by the smoothing can be quantified using Proposition 13.5 as we then have

$$f(\mathbf{u}) - \varepsilon \mathcal{U}_{\Omega} \leq f_{\Omega}(\mathbf{u}) \leq f(\mathbf{u}) - \varepsilon \mathcal{L}_{\Omega},$$

provided that $\mathcal{L}_{\Omega} \leq \Omega(\mathbf{v}) \leq \mathcal{U}_{\Omega}$ for all $\mathbf{v} \in \text{dom}(\Omega)$.

Primal approach

Following Definition 13.2, if we use dual regularization $\varepsilon\Omega$, where $\varepsilon > 0$ controls the regularization strength, the corresponding primal regularization is $R = \varepsilon\Omega^*(\cdot/\varepsilon)$. That is, we have

$$f_{\varepsilon\Omega} = f \square \varepsilon\Omega^*(\cdot/\varepsilon).$$

In the particular case $\Omega(\mathbf{v}) = \frac{1}{2}\|\mathbf{v}\|_2^2$, we have

$$R(\mathbf{u}) = \frac{\varepsilon}{2}\|\mathbf{u}/\varepsilon\|_2^2 = \frac{1}{2\varepsilon}\|\mathbf{u}\|_2^2 = \frac{1}{\varepsilon}\Omega(\mathbf{u}).$$

We therefore get

$$f_{\varepsilon\Omega} = f \square \frac{1}{\varepsilon}\Omega = \frac{1}{\varepsilon}(\varepsilon f \square \Omega) = \frac{1}{\varepsilon}\text{env}_{\varepsilon}f.$$

13.3.5 Generalized entropies

A natural choice of dual regularization $\Omega(\boldsymbol{\pi})$, when $\boldsymbol{\pi} \in \Delta^M$ is a discrete probability distribution, is a negative entropy function, also known as **negentropy**. Since negentropies play a major role in smoothed max operators, we discuss them in detail here.

Information content and entropy

An entropy function measures the amount of “surprise” of a random variable or equivalently of a distribution. To define an entropy, we must first define the **information content** $I(E)$ of an event E . The value returned by such a function should be 0 if the probability of the event is 1, as there is no surprise. Conversely, information content should attain its maximal value if the probability of the event is 0, as it is maximally surprising. Furthermore, the more probable an event E is, the less surprising it is. Therefore, when $p(E)$ increases, $I(E)$ should decrease. Overloading the notation, we also write the information content of the outcome y of a random variable Y as the information content of the event $\{Y = y\}$,

$$I(y) := I(\{Y = y\}).$$

Given an information content function, we can then define the **entropy** $H(Y)$ of a random variable $Y \in \mathcal{Y}$ as the expected information content,

$$H(Y) := \mathbb{E}[I(Y)].$$

Different definitions of information content lead to different definitions of entropy.

Shannon's entropy

A definition of information content satisfying the criteria above is

$$I(E) := \log \left(\frac{1}{p(E)} \right) = -\log p(E).$$

Indeed, $-\log 1 = 0$, $-\log 0 = \infty$ and $-\log$ is a decreasing function over $(0, 1]$. Using this information content definition leads to **Shannon's entropy** (Shannon, 1948)

$$H(Y) = \mathbb{E}[I(Y)] = - \sum_{y \in \mathcal{Y}} p(y) \log p(y).$$

We can therefore define the Shannon entropy of a discrete probability distribution $\pi \in \Delta^M$ as

$$H(\pi) = - \sum_{i=1}^M \pi_i \log \pi_i = -\langle \pi, \log \pi \rangle$$

and use the corresponding negentropy as regularization

$$\Omega(\pi) = -H(\pi) = \langle \pi, \log \pi \rangle.$$

The function is strongly convex w.r.t. $\|\cdot\|_1$ over Δ^M . However, it is not strongly convex over \mathbb{R}_+^M , since this is not a bounded set; see for instance (Blondel, 2019, Proposition 2). Since Ω is added to f^* in Eq. (13.3), we can therefore use this choice of Ω to smooth out a function f if $\text{dom}(f^*) \subseteq \Delta^M$.

Gini's entropy

As an alternative, we can define information content as

$$I(E) = \frac{1}{2}(1 - p(E)).$$

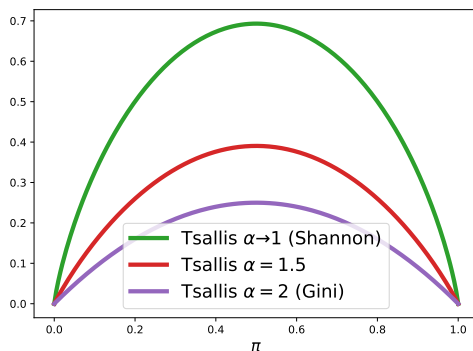


Figure 13.5: Tsallis entropies of the distribution $\pi = (\pi, 1 - \pi) \in \Delta^2$, for $\pi \in [0, 1]$. An entropy is a non-negative concave function that attains its maximum at the uniform distribution, here $(0.5, 0.5)$. A negative entropy, a.k.a. negentropy, can be used as a dual regularization function Ω to smooth out a function f when $\text{dom}(f^*) \subseteq \Delta^M$.

The $\frac{1}{2}$ factor is for later mathematical convenience. This again satisfies the criteria of an information content function. Indeed, i) when $p(E) = 1$, $I(E) = 0$ ii) when $p(E) = 0$, $I(E)$ attains its maximum of $\frac{1}{2}$ iii) the function is decreasing w.r.t. $p(E)$. Using this information content definition leads to **Gini's entropy** a.k.a. Gini index (Gini, 1912)

$$H(Y) = \mathbb{E}[I(Y)] = \frac{1}{2} \sum_{y \in \mathcal{Y}} p(y)(1 - p(y)).$$

We can use Gini's negative entropy to define for all $\pi \in \Delta^M$

$$\Omega(\pi) = \frac{1}{2} \langle \pi, \pi - \mathbf{1} \rangle = \frac{1}{2} (\|\pi\|_2^2 - 1).$$

The function is strongly convex w.r.t. $\|\cdot\|_2$ over \mathbb{R}^M . We can therefore use this choice of Ω to smooth out a function f if $\text{dom}(f^*) \subseteq \mathbb{R}^M$. This means that the set of functions that we can smooth out with Gini entropy is larger than the set of functions we can smooth out with Shannon entropy.

Tsallis entropies

Given $\alpha \geq 1$, a more general information content definition is

$$I(E) = \frac{1}{\alpha(\alpha - 1)} (1 - p(E)^{\alpha-1}).$$

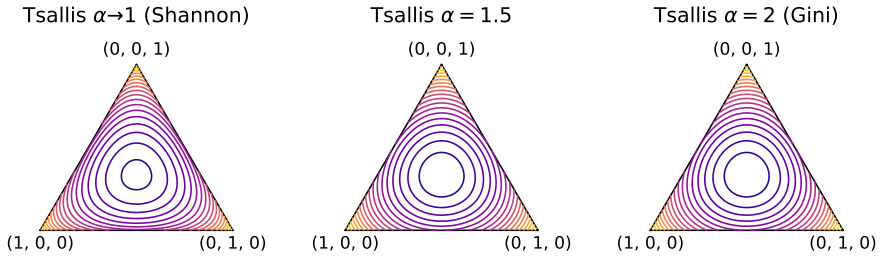


Figure 13.6: Contours of Tsallis entropies on the probability simplex.

Using this definition leads to the **Tsallis entropy** (Tsallis, 1988)

$$H(Y) = \mathbb{E}[I(Y)] = \frac{1}{\alpha(\alpha - 1)} \sum_{y \in \mathcal{Y}} p(y)(1 - p^{\alpha-1}(y)).$$

The Tsallis entropy recovers the Shannon entropy in the limit $\alpha \rightarrow 1$ and the Gini entropy when $\alpha = 2$. We can use the Tsallis negative entropy to define for all $\boldsymbol{\pi} \in \Delta^M$

$$\Omega(\boldsymbol{\pi}) = \frac{1}{\alpha(\alpha - 1)} \langle \boldsymbol{\pi}, \boldsymbol{\pi}^{\alpha-1} - \mathbf{1} \rangle = \frac{1}{\alpha(\alpha - 1)} (\|\boldsymbol{\pi}\|_\alpha^\alpha - 1),$$

where $\|\mathbf{v}\|_p$ is the p -norm for ($p \geq 1$)

$$\|\mathbf{v}\|_p := \left(\sum_{i=1}^M v_i^p \right)^{\frac{1}{p}},$$

so that

$$\|\mathbf{v}\|_p^p = \sum_{i=1}^M v_i^p.$$

Tsallis entropies for $\alpha \rightarrow 1$ (Shannon entropy), $\alpha = 1.5$ and $\alpha = 2$ (Gini entropy) are illustrated in Fig. 13.5 and Fig. 13.6.

Definition and properties of generalized entropies

So far, we saw how to define an entropy as the expected information content. However, generalized entropies (DeGroot, 1962; Grünwald and Dawid, 2004) do not necessarily need to take this form. We follow the definition of Blondel *et al.* (2020).

Definition 13.4 (Entropy function). A function $H: \Delta^M \rightarrow \mathbb{R}_+$ is an entropy if

1. $H(\boldsymbol{\pi}) = 0$ if $\boldsymbol{\pi} \in \{\mathbf{e}_1, \dots, \mathbf{e}_M\}$,
2. H is strictly concave,
3. $H(P\boldsymbol{\pi}) = H(\boldsymbol{\pi})$ for any permutation matrix P .

This definition implies that H is non-negative and is uniquely maximized by the uniform distribution (Blondel *et al.*, 2020, Proposition 4). This is indeed what we expect from an entropy function. An example is the squared p -norm entropy (Blondel *et al.*, 2020)

$$H(\boldsymbol{\pi}) = \frac{1}{2} - \frac{1}{2} \|\boldsymbol{\pi}\|_p^2.$$

Since the squared p -norm is strongly convex for $p \in (1, 2]$ (Ball *et al.*, 2002), this entropy is strongly concave for $p \in (1, 2]$ and can therefore be used to smooth out functions.

We now illustrate how to apply these techniques to compute smoothed ReLUs and smoothed max operators, as well as continuous relaxations of step functions and argmax operators.

13.4 Smoothed ReLU functions

To demonstrate the application of the smoothing techniques discussed in this chapter, we begin by explaining how to smooth the ReLU function. The ReLU function is defined by

$$\text{relu}(u) := \begin{cases} u & \text{if } u \geq 0 \\ 0 & \text{otherwise} \end{cases} = \max(u, 0).$$

We recall that in order to smooth a function f by the dual approach, we calculate its conjugate f^* , add regularization Ω to it to obtain $f_\Omega^* := f^* + \Omega$ and then obtain f_Ω by computing f_Ω^{**} .

Here, we wish to smooth out $f = \text{relu}$. Its convex conjugate is

$$\text{relu}^*(\pi) = \iota_{[0,1]}(\pi) = \begin{cases} 0 & \text{if } \pi \in [0, 1] \\ \infty & \text{if } \pi \notin [0, 1] \end{cases}.$$

To notice why, we observe that

$$\text{relu}(u) = \max_{\pi \in [0,1]} u \cdot \pi = \max_{\pi \in \{0,1\}} u \cdot \pi = \begin{cases} u & \text{if } u \geq 0 \\ 0 & \text{otherwise} \end{cases}. \quad (13.4)$$

Indeed, since the objective is linear in π , the maximum is attained at one of the extreme points of $[0, 1]$, so that we can replace the constraint $\pi \in [0, 1]$ with $\pi \in \{0, 1\}$. This shows that the ReLU is exactly the support function of $[0, 1]$. Since the conjugate of the support function is the indicator function, we indeed obtain $\text{relu}^* = \iota_{[0,1]}$. We therefore have

$$\text{relu}_\Omega^*(\pi) = \text{relu}^*(\pi) + \Omega(\pi) = \iota_{[0,1]}(\pi) + \Omega(\pi)$$

and for some choice of Ω , we need to be able to compute

$$\begin{aligned} \text{relu}_\Omega(u) &= \max_{\pi \in \mathbb{R}} u \cdot \pi - (\iota_{[0,1]}(\pi) + \Omega(\pi)) \\ &= \max_{\pi \in [0,1]} u \cdot \pi - \Omega(\pi). \end{aligned}$$

The softplus

If we use the regularizer $\Omega(\pi) = \pi \log \pi + (1 - \pi) \log(1 - \pi)$, which comes from using Shannon's negentropy $\langle \boldsymbol{\pi}, \log \boldsymbol{\pi} \rangle$ with $\boldsymbol{\pi} = (\pi, 1 - \pi)$, we obtain

$$\text{relu}_\Omega(u) = \text{softplus}(u) = \log(1 + \exp(u)).$$

This result is a special case of Proposition 13.9.

The sparseplus

If we use the regularizer $\Omega(\boldsymbol{\pi}) = \pi(\pi - 1)$, which comes from using Gini's negentropy with $\frac{1}{2} \langle \boldsymbol{\pi}, \boldsymbol{\pi} - \mathbf{1} \rangle$ with $\boldsymbol{\pi} = (\pi, 1 - \pi)$, we obtain

$$\text{relu}_\Omega(u) = \text{sparseplus}(u) = \begin{cases} 0, & u \leq -1 \\ \frac{1}{4}(u + 1)^2, & -1 < u < 1 \\ u, & u \geq 1 \end{cases}.$$

See Fig. 13.8 (left figure) for a comparison of softplus and sparseplus.

13.5 Smoothed max operators

As a more elaborate application of the smoothing techniques discussed in this chapter, we explain how to smooth max operators. Smoothed max operators include smoothed ReLU functions as a special case.

13.5.1 Definition and properties

With a slight notation overloading, given a vector $\mathbf{u} = (u_1, \dots, u_M) \in \mathbb{R}^M$, we define its maximum as

$$\max(\mathbf{u}) := \max_{j \in [M]} u_j.$$

To obtain a smooth approximation \max_Ω of \max , we again apply the dual approach. The conjugate of \max is

$$\max^*(\boldsymbol{\pi}) = \iota_{\Delta^M}(\boldsymbol{\pi}).$$

To notice why, we observe that the vertices of the probability simplex Δ^M are the standard basis vectors $\mathbf{e}_1, \dots, \mathbf{e}_M$. Since the objective is linear, we then have

$$\max(\mathbf{u}) = \max_{\boldsymbol{\pi} \in \Delta^M} \langle \mathbf{u}, \boldsymbol{\pi} \rangle = \max_{\boldsymbol{\pi} \in \{\mathbf{e}_1, \dots, \mathbf{e}_M\}} \langle \mathbf{u}, \boldsymbol{\pi} \rangle.$$

In other words, the maximum operator is exactly the support function of Δ^M . Since the conjugate of the support function is the indicator function, we indeed obtain $\max^* = \iota_{\Delta^M}$. We can therefore write

$$\max_\Omega^*(\boldsymbol{\pi}) = \max^*(\boldsymbol{\pi}) + \Omega(\boldsymbol{\pi}) = \Omega(\boldsymbol{\pi}) + \iota_{\Delta^M}(\boldsymbol{\pi})$$

and

$$\begin{aligned} \max_\Omega(\mathbf{u}) &= (\Omega + \iota_{\Delta^M})^*(\mathbf{u}) \\ &= \max_{\boldsymbol{\pi} \in \mathbb{R}^M} \langle \mathbf{u}, \boldsymbol{\pi} \rangle - (\Omega(\boldsymbol{\pi}) + \iota_{\Delta^M}(\boldsymbol{\pi})) \\ &= \max_{\boldsymbol{\pi} \in \Delta^M} \langle \mathbf{u}, \boldsymbol{\pi} \rangle - \Omega(\boldsymbol{\pi}). \end{aligned}$$

The smoothed max operator \max_Ω can be useful in a neural network, for example as a smoothed max pooling layer. Its properties have been studied in (Mensch and Blondel, 2018, Lemma 1), as we recall here for convenience.

Proposition 13.7 (Properties of \max_Ω). The following properties hold.

1. **Bounds:** if $\mathcal{L}_\Omega \leq \Omega(\boldsymbol{\pi}) \leq \mathcal{U}_\Omega$ for all $\boldsymbol{\pi} \in \Delta^M$, then $\max(\mathbf{u}) - \mathcal{U}_\Omega \leq \max_\Omega(\mathbf{u}) \leq \max(\mathbf{u}) - \mathcal{L}_\Omega$ for all $\mathbf{u} \in \mathbb{R}^M$.
2. **Monotonicity:** if $\mathbf{u} \leq \mathbf{v}$ (element-wise), then $\max_\Omega(\mathbf{u}) \leq \max_\Omega(\mathbf{v})$.
3. **Commutativity:** if $\Omega(P\boldsymbol{\pi}) = \Omega(\boldsymbol{\pi})$ for any permutation matrix P and any $\boldsymbol{\pi} \in \Delta^M$, then $\max_\Omega(P\mathbf{u}) = \max_\Omega(\mathbf{u})$ for any permutation matrix P .
4. **Distributivity of $+$:** $\max_\Omega(\mathbf{u} + c\mathbf{1}) = \max_\Omega(\mathbf{u}) + c$ for all $\mathbf{u} \in \mathbb{R}^M$ and all $c \in \mathbb{R}$.

These properties are leveraged in (Mensch and Blondel, 2018) to create differentiable dynamic programs. We consider in the following two possible choices of Ω leading to the softmax and sparsemax operators illustrated in Fig. 13.7.

Smoothed min operators

The minimum operator can be expressed in terms of the maximum operator, since for all $\mathbf{u} \in \mathbb{R}^M$,

$$\min(\mathbf{u}) = -\max(-\mathbf{u}).$$

Given a smoothed max operator \max_Ω , we can therefore easily define a smoothed min operator as

$$\min_\Omega(\mathbf{u}) := -\max_\Omega(-\mathbf{u}).$$

13.5.2 Reduction to root finding

Computing $\max_\Omega(\mathbf{u})$ for a general strongly-convex regularization Ω involves the resolution of a maximum over probability simplex constraints. For convenience, let us define the notation

$$\delta_\Omega(\mathbf{u}) := (\Omega + \iota_{\mathbb{R}_+^M})^*(\mathbf{u}) = \max_{\mathbf{v} \in \mathbb{R}_+^M} \langle \mathbf{u}, \mathbf{v} \rangle - \Omega(\mathbf{v}).$$

The following proposition shows that we can reduce computing \max_{Ω} to solving a root equation involving δ_{Ω} .

Proposition 13.8 (Computing \max_{Ω} as root finding). Suppose Ω is strongly convex. For all $\mathbf{u} \in \mathbb{R}^M$,

$$\begin{aligned}\max_{\Omega}(\mathbf{u}) &= \min_{\tau \in \mathbb{R}} \tau + \delta_{\Omega}(\mathbf{u} - \tau \mathbf{1}) \\ &= \tau^* + \delta_{\Omega}(\mathbf{u} - \tau^* \mathbf{1})\end{aligned}$$

and

$$\nabla \max_{\Omega}(\mathbf{u}) = \nabla \delta_{\Omega}(\mathbf{u} - \tau^* \mathbf{1}),$$

where τ^* is the solution w.r.t. τ of the above min, which satisfies the root equation

$$\langle \nabla \delta_{\Omega}(\mathbf{u} - \tau^* \mathbf{1}), \mathbf{1} \rangle = 1.$$

Proof. The idea is to keep the non-negativity constraint explicit, but to use a Lagrange multiplier for the equality constraint of the probability simplex. We then have

$$\begin{aligned}\max_{\Omega}(\mathbf{u}) &= \max_{\mathbf{v} \in \Delta^M} \langle \mathbf{u}, \mathbf{v} \rangle - \Omega(\mathbf{v}) \\ &= \max_{\mathbf{v} \in \mathbb{R}_+^M} \min_{\tau \in \mathbb{R}} \langle \mathbf{u}, \mathbf{v} \rangle - \Omega(\mathbf{v}) - \tau(\langle \mathbf{v}, \mathbf{1} \rangle - 1) \\ &= \min_{\tau \in \mathbb{R}} \tau + \max_{\mathbf{v} \in \mathbb{R}_+^M} \langle \mathbf{u} - \tau \mathbf{1}, \mathbf{v} \rangle - \Omega(\mathbf{v}) \\ &= \min_{\tau \in \mathbb{R}} \tau + \delta_{\Omega}(\mathbf{u} - \tau \mathbf{1}),\end{aligned}$$

where we used that we can swap the min and the max, since $(\mathbf{u}, \mathbf{v}) \mapsto \langle \mathbf{u}, \mathbf{v} \rangle - \Omega(\mathbf{v})$ is convex-concave and $\mathbf{v} \in \Delta^M$ is an affine constraint. The gradient $\nabla \delta_{\Omega}(\mathbf{u})$ follows from Danskin's theorem. The root equation follows from computing the derivative of $\tau \mapsto \tau + \delta_{\Omega}(\mathbf{u} - \tau \mathbf{1})$ and setting it to zero. \square

13.5.3 The softmax

When Ω is Shannon's negentropy, we obtain that \max_{Ω} is the softmax, already briefly discussed in Section 4.4.2.

Proposition 13.9 (Analytical expression of the softmax). When $\Omega(\boldsymbol{\pi}) = \langle \boldsymbol{\pi}, \log \boldsymbol{\pi} \rangle$, we get

$$\begin{aligned} \text{softmax}(\mathbf{u}) &:= \max_{\Omega}(\mathbf{u}) \\ &= \max_{\boldsymbol{\pi} \in \Delta^M} \langle \mathbf{u}, \boldsymbol{\pi} \rangle - \Omega(\boldsymbol{\pi}) \\ &= \text{logsumexp}(\mathbf{u}) \\ &= \log \sum_{j=1}^M e^{u_j}. \end{aligned}$$

Proof. Since $\text{dom}(\Omega) = \mathbb{R}_+^M$, we have $\delta_{\Omega} = \Omega^*$ (i.e., the non-negativity constraint is redundant). From Example 13.2, we therefore have $\delta_{\Omega}(\mathbf{u}) = \sum_{j=1}^M \exp(u_j - 1)$. From Proposition 13.8, $\max_{\Omega}(\mathbf{u}) = \tau^* + \delta_{\Omega}(\mathbf{u} - \tau^* \mathbf{1})$ where τ^* satisfies $\langle \nabla \delta_{\Omega}(\mathbf{u} - \tau^* \mathbf{1}), \mathbf{1} \rangle = 1$. Since $\nabla \delta_{\Omega}(\mathbf{u}) = \exp(\mathbf{u} - \mathbf{1})$, we need to solve $\sum_{j=1}^M \exp(u_j - 1 - \tau) = 1$. We therefore get $\tau^* + 1 = \text{logsumexp}(\mathbf{u})$ and therefore $\max_{\Omega}(\mathbf{u}) = \text{logsumexp}(\mathbf{u}) - 1 + \sum_{j=1}^M \exp(u_j - \text{logsumexp}(\mathbf{u})) = \text{logsumexp}(\mathbf{u})$. \square

Since $-\log M \leq \Omega(\boldsymbol{\pi}) \leq 0$ for all $\boldsymbol{\pi} \in \Delta^M$, following Proposition 13.7, we get for all $\mathbf{u} \in \mathbb{R}^M$

$$\max(\mathbf{u}) \leq \text{softmax}(\mathbf{u}) \leq \max(\mathbf{u}) + \log M.$$

A unique property of the softmax, which is not the case of all \max_{Ω} operators, is that it supports **associativity**.

Proposition 13.10 (Associativity of the softmax). For all $a, b, c \in \mathbb{R}$,

$$\text{softmax}(\text{softmax}(a, b), c) = \text{softmax}(a, \text{softmax}(b, c)).$$

13.5.4 The sparsemax

Alternatively, choosing Ω to be Gini's negentropy leads to the sparsemax (Martins and Astudillo, 2016; Mensch and Blondel, 2018).

Proposition 13.11 (Variational formulation of sparsemax). When $\Omega(\pi) = \frac{1}{2}\langle \pi, \pi - \mathbf{1} \rangle$, we have

$$\begin{aligned} \text{sparsemax}(\mathbf{u}) &:= \max_{\Omega}(\mathbf{u}) \\ &= \max_{\pi \in \Delta^M} \langle \mathbf{u}, \pi \rangle - \Omega(\pi) \\ &= \langle \mathbf{u}, \pi^* \rangle - \Omega(\pi^*) \end{aligned}$$

where

$$\pi^* = \text{sparseargmax}(\mathbf{u}) := \arg \min_{\pi \in \Delta^M} \|\mathbf{u} - \pi\|_2^2.$$

Proof. This follows from the fact that $\Omega(\pi)$ is up to a constant equal to $\frac{1}{2}\|\pi\|_2^2$ and completing the square. \square

Therefore, computing the sparsemax can use the sparseargmax (the Euclidean projection onto the probability simplex) as a building block. We discuss how to compute it in more detail in Section 13.7. Applying Proposition 13.8 gives an alternative formulation.

Proposition 13.12 (Sparsemax as root finding). When $\Omega(\pi) = \frac{1}{2}\langle \pi, \pi - \mathbf{1} \rangle$, we have

$$\text{sparsemax}(\mathbf{u}) = \max_{\Omega}(\mathbf{u}) = \min_{\tau \in \mathbb{R}} \tau + \frac{1}{2} \sum_{i=1}^M [u_i - \tau]_+^2$$

and τ^* satisfies

$$\sum_{i=1}^M [u_i - \tau]_+ = 1.$$

Proof. First, we compute the expression of $\delta_{\Omega}(\mathbf{u}) = \max_{\mathbf{v} \in \mathbb{R}_+^M} \langle \mathbf{u}, \mathbf{v} \rangle - \Omega(\mathbf{v})$. Setting the gradient of $\mathbf{v} \mapsto \langle \mathbf{u}, \mathbf{v} \rangle - \Omega(\mathbf{v})$ and clipping, we obtain $\mathbf{v}^* = [\mathbf{u}]_+$. Plugging \mathbf{v}^* back, we obtain $\delta_{\Omega}(\mathbf{u}) = \frac{1}{2} \sum_{i=1}^M [u_i]_+^2$. Using Proposition 13.8 proves the proposition's first part. Setting the derivative w.r.t. τ to zero gives the second part. \square

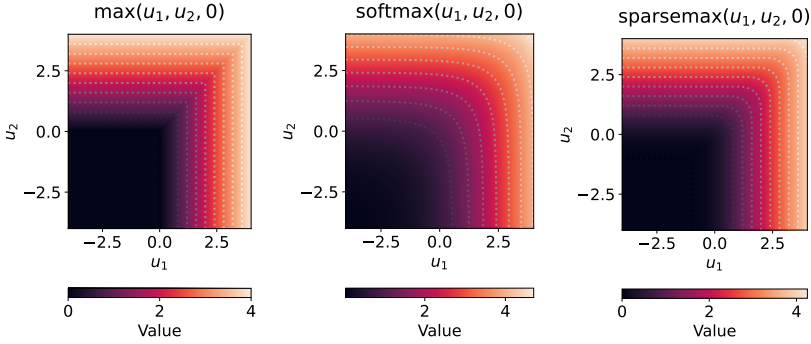


Figure 13.7: Max, softmax and sparsemax functions. The max function has non-smooth contour lines (set of points $\{\mathbf{u} \in \mathbb{R}^3 : f(\mathbf{u}) = c\}$ for some constant c represented by dashed gray lines). So the gradient along these contour lines switch suddenly at the corners of the contour lines switch. This shows that the max function is not differentiable everywhere, namely, non-differentiable on the set of points $\{\mathbf{u} \in \mathbb{R}^3 : u_i = u_j \text{ for any } i \neq j\}$. The contour lines of the softmax and sparsemax functions on the other hand are smooth illustrating that these functions are smooth counterpart of the max function.

It can be shown (Duchi *et al.*, 2008; Condat, 2016) that the exact solution τ^* is obtained by

$$\tau^* = \frac{1}{j^*} \left(\sum_{i=1}^{j^*} u_{[i]} - 1 \right), \quad (13.5)$$

where j^* is the largest $j \in [M]$ such that

$$u_j - \frac{1}{j} \left(\sum_{i=1}^j u_{[i]} - 1 \right) > 0,$$

and where we used the notation $u_{[1]} \geq u_{[2]} \geq \dots \geq u_{[M]}$. As an alternative, we can also compute τ^* approximately using a bisection or by gradient descent w.r.t. τ .

Since $\frac{1}{2M} \leq \|\boldsymbol{\pi}\|_2^2 \leq \frac{1}{2}$, we get $-\frac{M-1}{2M} \leq \|\boldsymbol{\pi}\|_2^2 \leq 0$ for all $\boldsymbol{\pi} \in \Delta^M$. Following Proposition 13.7, we therefore get for all $\mathbf{u} \in \mathbb{R}^M$

$$\max(\mathbf{u}) \leq \text{sparsemax}(\mathbf{u}) \leq \max(\mathbf{u}) + \frac{M-1}{2M}.$$

13.5.5 Recovering smoothed ReLU functions

Using the vector $\mathbf{u} = (u, 0) \in \mathbb{R}^2$ as input, the smoothed max operator recovers the smoothed ReLU:

$$\max_{\Omega}((u, 0)) = \text{relu}_{\Psi}(u),$$

where we defined $\Psi(\pi) := \Omega((\pi, 1 - \pi))$. With Ω being Shannon's negentropy, we recover $\Psi(\pi) = \pi \log \pi + (1 - \pi) \log(1 - \pi)$; with Ω being Gini's negentropy, we recover $\Psi(\pi) = \pi(\pi - 1)$, that we used to smooth the ReLU.

13.6 Relaxed step functions (sigmoid)

We now turn to creating continuous relaxations of step functions. The binary step function, a.k.a. Heaviside step function, is defined by

$$\text{step}(u) := \begin{cases} 1 & \text{if } u \geq 0 \\ 0 & \text{otherwise} \end{cases}.$$

From Eq. (13.4), its variational form is

$$\text{step}(u) = \arg \max_{\pi \in [0,1]} u \cdot \pi.$$

We can therefore define the relaxation

$$\text{step}_{\Omega}(u) := \arg \max_{\pi \in [0,1]} u \cdot \pi - \Omega(\pi).$$

Notice that, unlike the case of the smoothed ReLU, it is a regularized argmax, not a regularized max. Following Remark 13.1, strongly convex regularization Ω ensures that $\text{step}_{\Omega}(u)$ is a Lipschitz continuous function of u and is therefore, at least, differentiable almost everywhere, unlike $\text{step}(u)$.

The logistic function

If we use the regularizer $\Omega(\pi) = \pi \log \pi + (1 - \pi) \log(1 - \pi)$, we obtain the closed form

$$\text{step}_{\Omega}(u) = \text{logistic}(u) := \frac{1}{1 + e^{-u}} = \frac{e^u}{1 + e^u}.$$

This function is differentiable everywhere.

The sparse sigmoid

As an alternative, if we use $\Omega(\pi) = \pi(\pi - 1)$, we obtain a piecewise linear sigmoid,

$$\text{step}_\Omega(u) = \text{sparsesigmoid}(u) := \begin{cases} 0, & u \leq -1 \\ \frac{1}{2}(u + 1), & -1 < u < 1 \\ 1, & u \geq 1 \end{cases}.$$

Unlike the logistic function, it can reach the exact values 0 or 1. However, the function has two kinks, where the function is non-differentiable.

Link between smoothed ReLU functions and sigmoids

It turns out that the three sigmoids we presented above (step, logistic, sparsesigmoid) are all equal to the derivative of their corresponding smoothed ReLU function:

$$\begin{aligned} \text{step}(u) &= \text{relu}'(u) \\ \text{logistic}(u) &= \text{softplus}'(u) \\ \text{sparsesigmoid}(u) &= \text{sparseplus}'(u) \end{aligned}$$

and more generally

$$\text{relu}'_\Omega(u) = \text{step}_\Omega(u).$$

This is a consequence of Danskin's theorem; see Example 11.2. We illustrate the smoothed ReLU functions and relaxed step functions (sigmoids) in Fig. 13.8.

13.7 Relaxed argmax operators

We now turn to argmax operators, which are a generalization of step functions. With a slight notation overloading, let us now define

$$\text{argmax}(\mathbf{u}) := \phi(\arg \max_{j \in [M]} u_j),$$

where $\phi(j) = \text{onehot}(j) = \mathbf{e}_j$ is used to embed any integer $j \in [M]$ into \mathbb{R}^M . Following the previous discussion, we have the variational form

$$\text{argmax}(\mathbf{u}) = \arg \max_{\boldsymbol{\pi} \in \Delta^M} \langle \mathbf{u}, \boldsymbol{\pi} \rangle = \arg \max_{\boldsymbol{\pi} \in \{\mathbf{e}_1, \dots, \mathbf{e}_M\}} \langle \mathbf{u}, \boldsymbol{\pi} \rangle,$$

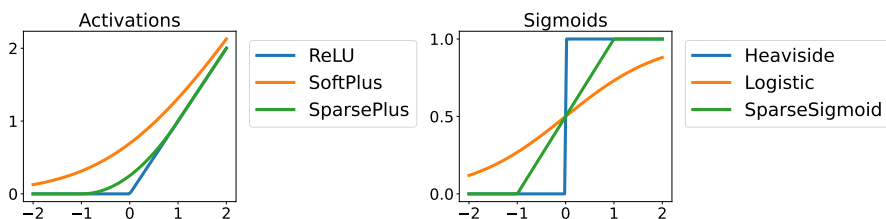


Figure 13.8: Smoothed ReLU functions and relaxed step functions (sigmoids). Differentiating the left functions gives the right functions.

where the second equality uses that a linear function is maximized at one of the vertices of the simplex. This variational form suggests to define the relaxation

$$\operatorname{argmax}_{\Omega}(\mathbf{u}) := \arg \max_{\boldsymbol{\pi} \in \Delta^M} \langle \mathbf{u}, \boldsymbol{\pi} \rangle - \Omega(\boldsymbol{\pi}).$$

Again, following Remark 13.1, $\operatorname{argmax}_{\Omega}(\mathbf{u})$ is guaranteed to be, at least, a differentiable almost everywhere function of \mathbf{u} if Ω is strongly convex.

Similarly to sigmoids, it turns out that these mappings are equal to the gradient of their corresponding smoothed max operator:

$$\operatorname{argmax}_{\Omega}(\mathbf{u}) = \nabla \max_{\Omega}(\mathbf{u}).$$

This is again a consequence of Danskin's theorem.

The softargmax

When using Shannon's entropy $\Omega(\boldsymbol{\pi}) = \langle \boldsymbol{\pi}, \log \boldsymbol{\pi} \rangle$, we obtain

$$\operatorname{argmax}_{\Omega}(\mathbf{u}) = \operatorname{softargmax}(\mathbf{u}) = \frac{\exp(\mathbf{u})}{\sum_{j=1}^M \exp(u_j)},$$

which is differentiable everywhere.

Proof. We know that $\max_{\Omega}(\mathbf{u}) = \log \operatorname{sumexp}(\mathbf{u})$ and that $\nabla \max_{\Omega}(\mathbf{u}) = \operatorname{argmax}_{\Omega}(\mathbf{u})$. Differentiating $\log \operatorname{sumexp}(\mathbf{u})$ gives $\operatorname{softargmax}(\mathbf{u})$. \square

The sparseargmax

When using Gini's entropy $\Omega(\boldsymbol{\pi}) = \frac{1}{2}\langle \boldsymbol{\pi}, \boldsymbol{\pi} - \mathbf{1} \rangle$, which is up to a constant equal to $\frac{1}{2}\|\boldsymbol{\pi}\|_2^2$, we obtain the sparseargmax (Martins and Astudillo, 2016)

$$\begin{aligned} \operatorname{argmax}_{\Omega}(\mathbf{u}) &= \operatorname{sparseargmax}(\mathbf{u}) \\ &:= \arg \max_{\boldsymbol{\pi} \in \Delta^M} \langle \mathbf{u}, \boldsymbol{\pi} \rangle - \frac{1}{2} \langle \boldsymbol{\pi}, \boldsymbol{\pi} - \mathbf{1} \rangle \\ &= \arg \max_{\boldsymbol{\pi} \in \Delta^M} \langle \mathbf{u}, \boldsymbol{\pi} \rangle - \frac{1}{2} \|\boldsymbol{\pi}\|_2^2 \\ &= \arg \min_{\boldsymbol{\pi} \in \Delta^M} \|\mathbf{u} - \boldsymbol{\pi}\|_2^2, \end{aligned}$$

which is nothing but the Euclidean projection onto the probability simplex (see also Section 16.3). The Euclidean projection onto the probability simplex Δ^M can be computed exactly using a median-finding-like algorithm. The complexity is $O(M)$ expected time and $O(M \log M)$ worst-case time (Brucker, 1984; Michelot, 1986; Duchi *et al.*, 2008; Condat, 2016). Computing the Euclidean projection onto the probability simplex boils down to computing τ^* given in Eq. (13.5). Once we computed it, we have

$$\operatorname{sparseargmax}(\mathbf{u}) = [\mathbf{u} - \tau^*]_+,$$

At its name indicates, and as the above equation shows, sparseargmax is **sparse**, but it is only differentiable almost everywhere. Note that the operator is originally known as sparsemax (Martins and Astudillo, 2016), but this is a misnomer, as it is really an approximation of the argmax. Therefore, in analogy with the softmax, we use the name sparseargmax. We compare the argmax, softmax and sparseargmax in Fig. 13.9 and Fig. 13.10.

Relaxed argmin operators

The argmin operator can be expressed in terms of the argmax operator,

$$\arg \min(\mathbf{u}) = \arg \max(-\mathbf{u}).$$

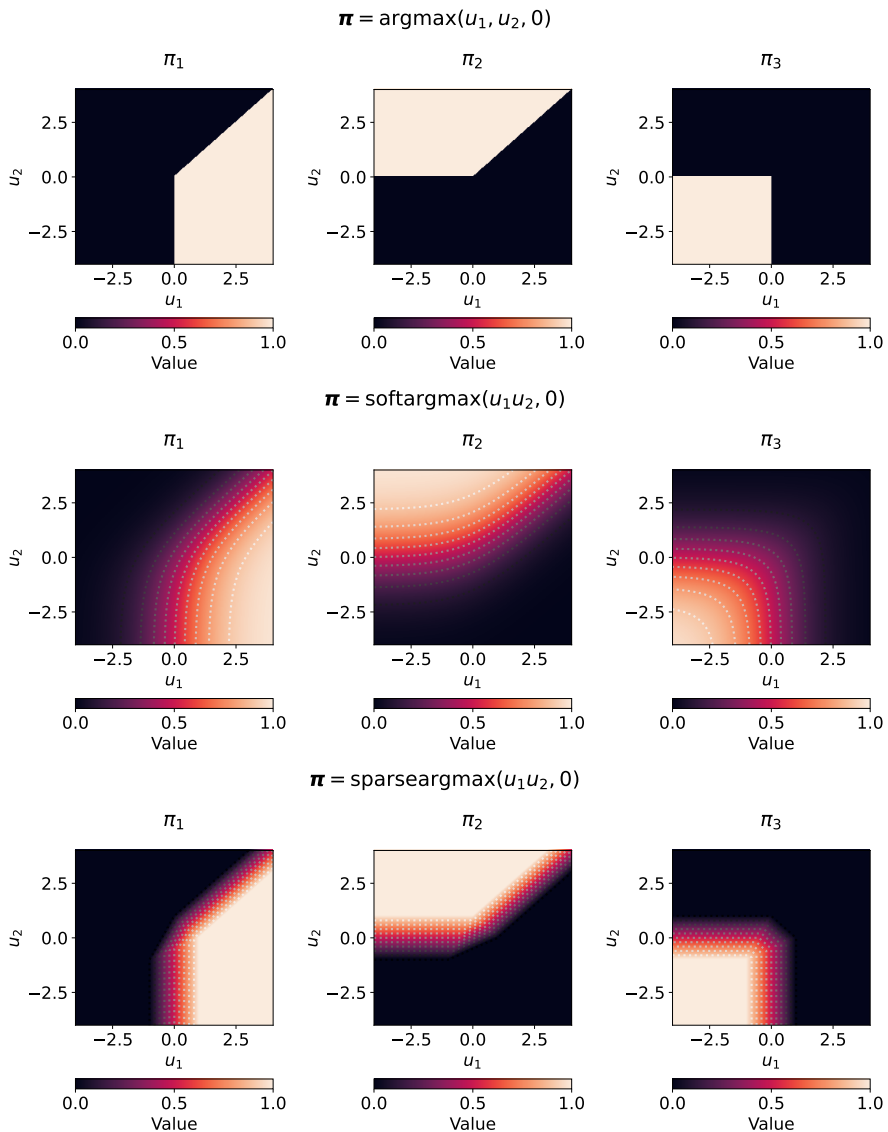


Figure 13.9: Values of $\operatorname{argmax}(\mathbf{u})$, $\operatorname{softargmax}(\mathbf{u})$, and $\operatorname{sparseargmax}(\mathbf{u})$ for $\mathbf{u} = (u_1, u_2, 0)$, when varying u_1 and u_2 . The argmax is a piecewise constant, discontinuous function. The $\operatorname{softargmax}$ is a continuous and differentiable everywhere function, but it is always strictly positive and therefore dense. The $\operatorname{sparseargmax}$ is a continuous function and its output can be sparse, but it is only a differentiable almost everywhere function.

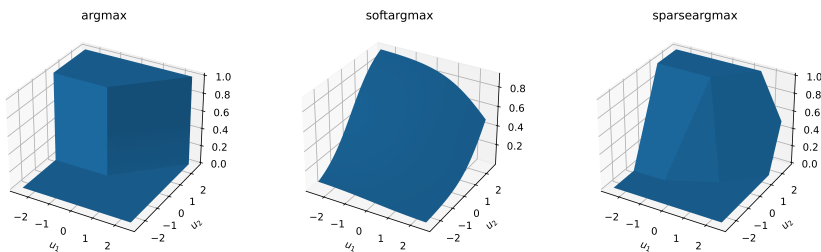


Figure 13.10: Same as Fig. 13.9 but using a 3D plot.

Given a relaxed argmax operator argmax_Ω , we can therefore define a relaxed argmin by

$$\text{argmin}_\Omega(\mathbf{u}) := \text{argmax}_\Omega(-\mathbf{u}).$$

We then have for all $\mathbf{u} \in \mathbb{R}^M$

$$\text{argmin}_\Omega(\mathbf{u}) = \nabla \text{min}_\Omega(\mathbf{u}).$$

13.8 Summary

- When a function f is non-differentiable (or worse, discontinuous), a reasonable approach is to replace it by its smooth approximation (or continuous relaxation).
- The first approach we reviewed is infimal convolution between f and primal regularization R . The Moreau envelope is a special case, obtained by using $R = \frac{1}{2} \|\cdot\|_2^2$.
- The second approach we reviewed is regularizing the convex conjugate f^* of f with some dual regularization Ω . We saw that the primal and dual approaches are equivalent when $R = \Omega^*$.
- The Legendre-Fenchel transformation, a.k.a. convex conjugate, can be seen as a dual representation of a function: instead of representing f by its graph $(\mathbf{u}, f(\mathbf{u}))$ for $\mathbf{u} \in \text{dom}(f)$, we can represent it by the set of tangents with slope \mathbf{v} and intercept

$-f^*(\mathbf{v})$ for $\mathbf{v} \in \text{dom}(f^*)$ As its name indicates, it is convex, even if the original function is not.

- We showed how to apply smoothing techniques to create smoothed ReLU functions and smoothed max operators. We also showed that taking their gradients allowed us to obtain generalized sigmoid functions and argmax operators.

14

Smoothing by integration

In this chapter, we review smoothing techniques based on **convolution**.

14.1 Convolution

14.1.1 Convolution operators

The convolution between two functions f and g produces another function, denoted $f * g$. It is defined by

$$(f * g)(\mu) := \int_{-\infty}^{\infty} f(u)g(\mu - u) du, \quad (14.1)$$

assuming that the integral is well defined. It is therefore the integral of the product of f and g after g is reflected about the y -axis and shifted. It can be seen as a generalization of the **moving average**. Using the change of variable $u := \mu + z$, which is again the **location-scale transform**, we can also write

$$(f * g)(\mu) = \int_{-\infty}^{\infty} f(\mu - z)g(z) dz = (g * f)(\mu). \quad (14.2)$$

The convolution operator is therefore **commutative**.

14.1.2 Convolution with a kernel

The convolution is frequently used together with a **kernel** κ to create a smooth approximation $f * \kappa$ of f . The most frequently used kernel is the **Gaussian kernel** with width σ , defined by

$$\kappa_\sigma(z) := \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{z}{\sigma}\right)^2}.$$

This is the probability density function (PDF) of the normal distribution with zero mean and variance σ^2 . The term $\frac{1}{\sqrt{2\pi}\sigma}$ is a normalization constant, ensuring that the kernel sums to 1 for all σ . We therefore say that κ_σ is a **normalized kernel**.

Averaging perspective

Applying the definition of the convolution in Eq. (14.1), we obtain

$$\begin{aligned} (f * \kappa_\sigma)(\mu) &:= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} f(u) e^{-\frac{1}{2}\left(\frac{\mu-u}{\sigma}\right)^2} du \\ &= \mathbb{E}_{U \sim p_{\mu,\sigma}}[f(U)], \end{aligned}$$

where

$$p_{\mu,\sigma}(u) := \kappa_\sigma(\mu - u) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{\mu-u}{\sigma}\right)^2}$$

is the PDF of the Gaussian distribution with mean μ and variance σ^2 . Therefore, we can see $f * \kappa_\sigma$ as the **expectation** of $f(u)$ over a Gaussian **centered around** μ . This property is true for all translation-invariant kernels, that correspond to a location-scale family distribution (e.g., the Laplace distribution). The convolution therefore performs an averaging with all points, with points nearby μ given more weight by the distribution. The parameter σ controls the importance we want to give to farther points. We call this viewpoint averaging, as we replace $f(\mathbf{u})$ by $\mathbb{E}[f(U)]$.

Perturbation perspective

Conversely, using the alternative definition of the convolution operator in Eq. (14.2), which stems from the commutativity of the convolution,

we have

$$\begin{aligned}(f * \kappa_\sigma)(\mu) &:= \int_{-\infty}^{\infty} f(\mu - z) e^{-\frac{1}{2}(\frac{z}{\sigma})^2} dz \\ &= \mathbb{E}_{Z \sim p_{0,\sigma}}[f(\mu - Z)] \\ &= \mathbb{E}_{Z \sim p_{0,\sigma}}[f(\mu + Z)],\end{aligned}$$

where, in the third line, we used that $p_{0,\sigma}$ is sign invariant, i.e., $p_{0,\sigma}(z) = p_{0,\sigma}(-z)$. This viewpoint shows that smoothing by convolution with a Gaussian kernel can also be seen as injecting Gaussian **noise** or **perturbations** to the function's input.

Limit case

When $\sigma \rightarrow 0$, the kernel κ_σ converges to a Dirac delta function,

$$\lim_{\sigma \rightarrow 0} \kappa_\sigma(z) = \delta(z).$$

Since the Dirac delta is the multiplicative identity of the convolution algebra (this is also known as the sifting property), when $\sigma \rightarrow 0$, $f * \kappa_\sigma$ converges to f , i.e.,

$$\lim_{\sigma \rightarrow 0} (f * \kappa_\sigma)(u) = f(u).$$

14.1.3 Discrete convolution

Many times, we work with functions whose convolution does not have an analytical form. In these cases, we can use a discrete convolution on a grid of values. For two functions f and g defined over \mathbb{Z} , the discrete convolution is defined by

$$(f * g)[i] := \sum_{j=-\infty}^{\infty} f[j]g[i-j].$$

As for its continuous counterpart, the discrete convolution is commutative, namely,

$$(f * g)[i] = \sum_{j=-\infty}^{\infty} f[i-j]g[j] = (g * f)[i].$$

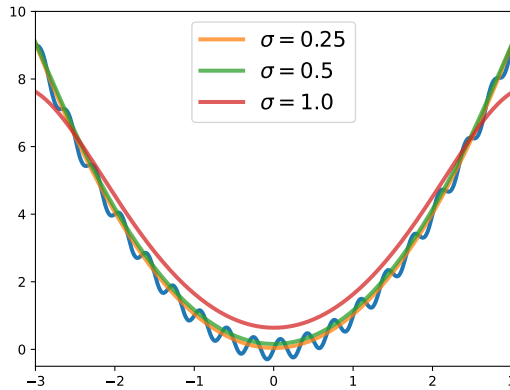


Figure 14.1: Smoothing of the signal $f[t] := t^2 + 0.3 \sin(6\pi t)$ with a sampled and renormalized Gaussian kernel.

When g has finite support over the set $S := \{-M, -M+1, \dots, 0, \dots, M-1, M\}$, meaning that $g[i] = 0$ for all $i \notin S$, a finite summation may be used instead, i.e.,

$$(f * g)[i] = \sum_{j=-M}^M f[i-j]g[j] = (g * f)[i].$$

In practice, convolution between a discrete signal $f: \mathbb{Z} \rightarrow \mathbb{R}$ and a continuous kernel $\kappa: \mathbb{R} \rightarrow \mathbb{R}$ is implemented by discretizing the kernel. One of the simplest approaches consists in sampling points on an interval, evaluating the kernel at these points and renormalizing the obtained values, so that the sampled kernel sums to 1. This is illustrated with the Gaussian kernel in Fig. 14.1. Since the Gaussian kernel decays exponentially fast, we can choose a small interval around 0. For a survey of other possible discretizations of the Gaussian kernel, see Getreuer (2013).

14.1.4 Differentiation

Remarkably, provided that the two functions are integrable with integrable derivatives, the derivative of the convolution satisfies

$$(f * g)' = (f' * g) = (f * g'),$$

which simply stems from switching derivative and integral in the definition of the convolution. Moreover, we have the following proposition.

Proposition 14.1 (Differentiability of the convolution). If g is n -times differentiable with compact support over \mathbb{R} and f is locally integrable over \mathbb{R} , then $f * g$ is n -times differentiable over \mathbb{R} .

14.1.5 Multidimensional convolution

So far, we studied the convolution of one-dimensional functions. The definition can be naturally extended to multidimensional functions $f: \mathbb{R}^M \rightarrow \mathbb{R}$ and $g: \mathbb{R}^M \rightarrow \mathbb{R}$ as

$$(f * g)(\boldsymbol{\mu}) := \int_{\mathbb{R}^M} f(\mathbf{u})g(\boldsymbol{\mu} - \mathbf{u}) d\mathbf{u},$$

assuming again that the integral exists. Typically, a Gaussian kernel with diagonal covariance matrix is used

$$\kappa_{\sigma}(\mathbf{z}) := \prod_{j=1}^M \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{1}{2}(\frac{z_j}{\sigma_j})^2} = \frac{1}{\sqrt{2\pi}^M \sigma^M} e^{-\frac{1}{2} \frac{\|\mathbf{z}\|_2^2}{\sigma^2}}, \quad (14.3)$$

where, in the second equality, we assumed $\sigma_1 = \dots = \sigma_M$. In an image processing context, where $M = 2$, it is approximated using a discrete convolution and it is called a **Gaussian blur**.

14.1.6 Link between convolution and infimal convolution

The infimal convolution we studied in Section 13.1 takes the form

$$(f \square g)(\boldsymbol{\mu}) := \inf_{\mathbf{u} \in \mathbb{R}^M} f(\mathbf{u}) + g(\boldsymbol{\mu} - \mathbf{u}).$$

In comparison, the classical convolution takes the form

$$(F * G)(\boldsymbol{\mu}) := \int_{\mathbb{R}^M} F(\mathbf{u})G(\boldsymbol{\mu} - \mathbf{u}) d\mathbf{u}.$$

The two forms of convolution are clearly related. Infimal convolution performs an infimum and uses the sum of f and g : it uses a **min-plus algebra**. Classical convolution performs an integral and uses the product of F and G : it uses a **sum-product algebra**.

14.1.7 The soft infimal convolution

The link between the infimal convolution and the classical convolution can be further elucidated if we replace the infimum with a soft minimum in the definition of the infimal convolution.

Definition 14.1 (Soft infimal convolution). The soft infimal convolution between $f: \mathbb{R}^M \rightarrow \mathbb{R}$ and $g: \mathbb{R}^M \rightarrow \mathbb{R}$ is

$$(f \square_{\varepsilon} g)(\boldsymbol{\mu}) := \operatorname{softmin}_{\varepsilon} \limits_{\mathbf{u} \in \mathbb{R}^M} f(\mathbf{u}) + g(\boldsymbol{\mu} - \mathbf{u}),$$

where we defined the soft minimum (assuming that it exists) over \mathcal{S} of any function $h: \mathcal{S} \rightarrow \mathbb{R}$ as

$$\operatorname{softmin}_{\varepsilon} \limits_{\mathbf{u} \in \mathcal{S}} h(\mathbf{u}) := -\varepsilon \log \int_{\mathcal{S}} \exp(-h(\mathbf{u})/\varepsilon) d\mathbf{u}.$$

We recover the infimal convolution as $\varepsilon \rightarrow 0$.

Computation using a convolution

We now show that we can rewrite the soft infimal convolution using a classical convolution. Indeed, by using the exponential change of variable (sometimes referred to as **Cole-Hopf transformation** in a partial differential equation context)

$$\begin{aligned} \mathcal{C}_{\varepsilon}\{f\}(\mathbf{u}) &:= \exp(-f(\mathbf{u})/\varepsilon) \\ \mathcal{C}_{\varepsilon}^{-1}\{F\}(\mathbf{v}) &= -\varepsilon \log F(\mathbf{v}), \end{aligned}$$

we can define each function in the exponential domain,

$$\begin{aligned} F_{\varepsilon} &:= \mathcal{C}_{\varepsilon}\{f\} \\ G_{\varepsilon} &:= \mathcal{C}_{\varepsilon}\{g\} \\ H_{\varepsilon} &:= \mathcal{C}_{\varepsilon}\{h_{\varepsilon}\}. \end{aligned}$$

It is easy to check that we then have

$$H_\varepsilon(\boldsymbol{\mu}) = (F_\varepsilon * G_\varepsilon)(\boldsymbol{\mu}).$$

Back to log domain, we obtain

$$h_\varepsilon(\boldsymbol{\mu}) = \mathcal{C}_\varepsilon^{-1}\{H_\varepsilon\}(\boldsymbol{\mu}).$$

Combining the transformation and its inverse, we can write

$$h_\varepsilon(\boldsymbol{\mu}) = \mathcal{C}_\varepsilon^{-1}\{\mathcal{C}_\varepsilon\{f\} * \mathcal{C}_\varepsilon\{g\}\}(\boldsymbol{\mu}).$$

What we have shown is that, after an exponential change of variable, the soft infimal convolution can be reduced to the computation of a convolution. This is useful as a discrete convolution on a grid of size n can be computed in $O(n \log n)$.

14.1.8 The soft Moreau envelope

We saw in Section 13.1.2 that the infimal convolution between f and $R(z) = \frac{1}{2}z^2$ is the Moreau envelope,

$$M_f(\boldsymbol{\mu}) := (f \square R)(\boldsymbol{\mu}) = \inf_{\mathbf{u} \in \mathbb{R}^M} f(\mathbf{u}) + \frac{1}{2}\|\boldsymbol{\mu} - \mathbf{u}\|_2^2.$$

Replacing the infimal convolution with a soft infimal convolution, we can define the “soft” Moreau envelope,

$$M_f^\varepsilon(\boldsymbol{\mu}) := (f \square_\varepsilon R)(\boldsymbol{\mu}) = \text{softmin}_\varepsilon f(\mathbf{u}) + \frac{1}{2}\|\boldsymbol{\mu} - \mathbf{u}\|_2^2.$$

We emphasize that this operation is **not** the same as the convolution of f with a Gaussian kernel. Indeed, we have

$$M_f^\varepsilon(\boldsymbol{\mu}) = -\varepsilon \log \int_{\mathbb{R}^M} \exp\left(\left(-f(\mathbf{u}) - \frac{1}{2}\|\boldsymbol{\mu} - \mathbf{u}\|_2^2\right)/\varepsilon\right) d\mathbf{u}.$$

while

$$(f * \kappa_\sigma)(\boldsymbol{\mu}) := \int_{\mathbb{R}^M} f(\mathbf{u}) \kappa_\sigma(\boldsymbol{\mu} - \mathbf{u}) d\mathbf{u},$$

where κ_σ is for instance defined in Eq. (14.3).

We saw that the Moreau envelope is a smooth function. One may therefore ask what do we gain from using a soft Moreau envelope. The benefit can be computational, as the latter can be approximated using a discrete convolution.

14.2 Fourier and Laplace transforms

Let us define the **Fourier transform** of f by

$$F(s) := \mathcal{F}\{f\}(s) := \int_{-\infty}^{\infty} f(t)e^{-i2\pi st} dt, \quad s \in \mathbb{R}.$$

Note that $\mathcal{F}\{f\}$ is a function transformation: it transforms f into another function F .

14.2.1 Convolution theorem

Now, consider the convolution

$$h(t) := (f * g)(t).$$

If we define the three transformations

$$F := \mathcal{F}\{f\}, \quad G := \mathcal{F}\{g\}, \quad H := \mathcal{F}\{h\},$$

the **convolution theorem** states that

$$H(s) = F(s) \cdot G(s), \quad s \in \mathbb{R}.$$

Written differently, we have

$$\mathcal{F}\{f * g\} = \mathcal{F}\{f\} \cdot \mathcal{F}\{g\}.$$

In words, in the Fourier domain, the convolution operation becomes a multiplication. Conversely,

$$h(t) = (f * g)(t) = \mathcal{F}^{-1}\{F \cdot G\}(t), \quad t \in \mathbb{R}.$$

The convolution theorem also holds if we replace the Fourier transform with the Laplace transform or with the two-sided (bilateral) Laplace transform.

14.2.2 Link between Fourier and Legendre transforms

In Section 13.2, we studied another function transformation: the convex conjugate, also known as Legendre-Fenchel transform. We recap the analogies between these transforms in Table 14.1. In particular, the counterpart of

$$\mathcal{F}\{f * g\} = \mathcal{F}\{f\} \cdot \mathcal{F}\{g\}.$$

Table 14.1: Analogy between Fourier and Legendre transforms. See Proposition 13.3 for more conjugate calculus rules.

	Fourier $\mathcal{F}\{f\}$	Legendre f^*
Semiring	$(+, \cdot)$	$(\min, +)$
Scaling ($a > 0$)	$f(t) = g(t/a)$ $\mathcal{F}\{f\}(t) = a\mathcal{F}\{g\}(as)$	$f(t) = ag(t/a)$ $f^*(s) = ag^*(s)$
Translation	$f(t) = g(t - t_0)$ $\mathcal{F}\{f\}(s) = e^{-i2\pi t_0 s} \mathcal{F}\{g\}(s)$	$f(t) = g(t - t_0)$ $f^*(s) = g^*(s) + t_0$
Convolution	$h = f * g$ $\mathcal{F}\{h\} = \mathcal{F}\{f\} \cdot \mathcal{F}\{g\}$	$h = f \square g$ $h^* = f^* + g^*$
Gaussian / quadratic	$f(t) = e^{-at^2}$ $\mathcal{F}\{f\}(s) = \sqrt{\frac{\pi}{a}} e^{-\pi^2 s^2 / a}$	$f(t) = \frac{a}{2} t^2$ $f^*(s) = \frac{1}{2a} s^2$
Smoothing	$f * \kappa_\sigma$	$f \square \frac{1}{2\varepsilon} \ \cdot\ _2^2$

for the infimal convolution is

$$(f \square g)^* = f^* + g^*.$$

In words, the Legendre-Fenchel transform is to the infimal convolution what the Fourier transform is to the convolution.

14.2.3 The soft Legendre-Fenchel transform

We saw in Section 13.2 that the Legendre-Fenchel transform (convex conjugate) of a function $f: \mathbb{R}^M \rightarrow \mathbb{R}$ is

$$f^*(\mathbf{v}) := \max_{\mathbf{u} \in \mathbb{R}^M} \langle \mathbf{u}, \mathbf{v} \rangle - f(\mathbf{u}).$$

If necessary, we can support constraints by including an indicator function in the definition of f . The conjugate can be smoothed out using a log-sum-exp, which plays the role of a soft maximum (Section 13.5).

Definition 14.2 (Soft convex conjugate).

$$f_\varepsilon^*(\mathbf{v}) := \operatorname{softmax}_\varepsilon \langle \mathbf{u}, \mathbf{v} \rangle - f(\mathbf{u}),$$

$$\mathbf{u} \in \mathbb{R}^M$$

where we defined the soft maximum (assuming that it exists) over \mathcal{S} of any function $g: \mathcal{S} \rightarrow \mathbb{R}$ as

$$\operatorname{softmax}_\varepsilon g(\mathbf{u}) := \varepsilon \log \int_{\mathcal{S}} \exp(g(\mathbf{u})/\varepsilon) d\mathbf{u}.$$

In the limit $\varepsilon \rightarrow 0$, we recover the convex conjugate.

Computation using a convolution

We now show that this smoothed conjugate can be rewritten using a convolution if we apply a bijective transformation to f .

Proposition 14.2 (Smoothed convex conjugate as convolution). The smoothed conjugate can be rewritten as

$$f_\varepsilon^*(\mathbf{v}) = \mathcal{Q}_\varepsilon^{-1} \left\{ \frac{1}{\mathcal{Q}_\varepsilon\{f\} * G_\varepsilon} \right\}(\mathbf{v})$$

where

$$G_\varepsilon := \mathcal{C}_\varepsilon \left\{ \frac{1}{2} \|\cdot\|_2^2 \right\} = \exp \left(-\frac{1}{2} \frac{\|\cdot\|_2^2}{\varepsilon} \right)$$

$$\mathcal{Q}_\varepsilon\{f\} := \mathcal{C}_\varepsilon \left\{ f(\cdot) - \frac{1}{2} \|\cdot\|_2^2 \right\} = \exp \left(\frac{1}{2\varepsilon} \|\cdot\|_2^2 - \frac{1}{\varepsilon} f(\cdot) \right)$$

$$\mathcal{Q}_\varepsilon^{-1}\{F\} := \frac{1}{2} \|\cdot\|_2^2 - \varepsilon \log(F(\cdot)).$$

This insight was [tweeted](#) by Gabriel Peyré in April 2020.

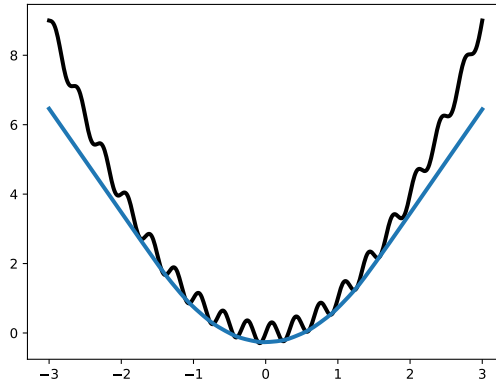


Figure 14.2: Applying the smoothed conjugate twice gives a smoothed biconjugate (convex envelope) of the function.

Proof.

$$\begin{aligned}
 f_\varepsilon(\mathbf{v}) &:= \varepsilon \log \int \exp \left(\frac{1}{\varepsilon} \langle \mathbf{u}, \mathbf{v} \rangle - \frac{1}{\varepsilon} f(\mathbf{u}) \right) d\mathbf{u} \\
 &= \varepsilon \log \int \exp \left(-\frac{1}{2\varepsilon} \|\mathbf{u} - \mathbf{v}\|_2^2 + \frac{1}{2\varepsilon} \|\mathbf{u}\|_2^2 + \frac{1}{2\varepsilon} \|\mathbf{v}\|_2^2 - \frac{1}{\varepsilon} f(\mathbf{u}) \right) d\mathbf{u} \\
 &= \varepsilon \log \int \exp \left(-\frac{1}{2\varepsilon} \|\mathbf{u} - \mathbf{v}\|_2^2 + \frac{1}{2\varepsilon} \|\mathbf{u}\|_2^2 - \frac{1}{\varepsilon} f(\mathbf{u}) \right) d\mathbf{u} + \frac{1}{2} \|\mathbf{v}\|_2^2 \\
 &= \varepsilon \log \int G_\varepsilon(\mathbf{v} - \mathbf{u}) \mathcal{Q}_\varepsilon\{f\}(\mathbf{u}) d\mathbf{u} + \frac{1}{2} \|\mathbf{v}\|_2^2 \\
 &= \varepsilon \log(\mathcal{Q}_\varepsilon\{f\} * G_\varepsilon)(\mathbf{v}) + \frac{1}{2} \|\mathbf{v}\|_2^2 \\
 &= \frac{1}{2} \|\mathbf{v}\|_2^2 - \varepsilon \log \left(\frac{1}{\mathcal{Q}_\varepsilon\{f\} * G_\varepsilon} \right) (\mathbf{v}) \\
 &= \mathcal{Q}_\varepsilon^{-1} \left\{ \frac{1}{\mathcal{Q}_\varepsilon\{f\} * G_\varepsilon} \right\} (\mathbf{v})
 \end{aligned}$$

□

What did we gain from this viewpoint? The convex conjugate can often be difficult to compute in closed form. If we replace \mathbb{R}^M with a discrete set \mathcal{S} (i.e., a grid), we can then approximate the smoothed convex

conjugate in $O(n \log n)$, where $n = |\mathcal{S}|$, using a discrete convolution,

$$\begin{aligned} (\mathcal{Q}_\varepsilon\{f\} * G_\varepsilon)(\mathbf{v}) &\approx \sum_{\mathbf{u} \in \mathcal{S}} G_\varepsilon(\mathbf{v} - \mathbf{u}) \mathcal{Q}_\varepsilon\{f\}(\mathbf{u}) \\ &= \mathbf{K} \mathbf{q}, \end{aligned}$$

where \mathbf{K} is the $n \times n$ Gaussian kernel matrix whose entries correspond to $\exp(-\frac{1}{2\varepsilon} \|\mathbf{u} - \mathbf{u}'\|_2^2)$ for $\mathbf{u}, \mathbf{u}' \in \mathcal{S}$ and \mathbf{q} is the n -dimensional vector whose entries correspond to $\exp(\frac{1}{\varepsilon}(\frac{1}{2}\|\mathbf{v}\|_2^2 - f(\mathbf{u})))$ for $\mathbf{u} \in \mathcal{S}$. This provides a GPU-friendly alternative to the fast Legendre transform algorithm, discussed in Section 13.2. Of course, due to the curse of dimensionality, the technique is limited to functions defined on low-dimensional sets. We illustrate in Fig. 14.2 the application of the technique to computing an approximate biconjugate (convex envelope) of a function.

Remark 14.1 (Link with the two-sided Laplace transform). For one-dimensional functions, instead of using a convolution, we can also write the soft convex conjugate as

$$\begin{aligned} f_\varepsilon^*(v) &= \varepsilon \log \int_{-\infty}^{\infty} \exp\left(\frac{1}{\varepsilon} [uv - f(u)]\right) du \\ &= \varepsilon \log \mathcal{B}\left\{e^{-\frac{f}{\varepsilon}}\right\}\left(-\frac{v}{\varepsilon}\right) \\ &= -\mathcal{C}_\varepsilon^{-1}\{\mathcal{B}\{\mathcal{C}_\varepsilon\{f\}\}\}\left(-\frac{v}{\varepsilon}\right) \end{aligned}$$

where we defined the two-sided (bilateral) Laplace transform

$$\mathcal{B}\{g\}(v) := \int_{-\infty}^{\infty} e^{-uv} g(u) du$$

and where we assumed that the integral exists.

14.3 Examples

In this section, we review practical examples for which the convolution with a Gaussian kernel enjoys an analytical solution.

14.3.1 Smoothed step function

Example 14.1 (Smoothed Heaviside). The Heaviside step function is defined by

$$\text{step}(u) := h(u) := \begin{cases} 1 & \text{if } u \geq 0 \\ 0 & \text{otherwise} \end{cases}.$$

With the Gaussian kernel, we therefore obtain

$$\begin{aligned} (h * \kappa_\sigma)(\mu) &= \int_{-\infty}^{\mu} \kappa_\sigma(z) h(\mu - z) dz + \int_{\mu}^{\infty} \kappa_\sigma(z) h(\mu - z) dz \\ &= \int_{-\infty}^{\mu} \kappa_\sigma(z) dz \\ &= \Phi_\sigma(\mu) \\ &= \frac{1}{2} \left[1 + \text{erf} \left(\frac{\mu}{\sqrt{2}\sigma} \right) \right], \end{aligned}$$

where $\Phi_\sigma(\mu)$ is the CDF of the Gaussian distribution with zero mean and variance σ^2 , and where we used the error function

$$\text{erf}(z) := \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt,$$

that we both already encountered in Chapter 3. Although there is no closed form for the error function, it is commonly available in numerical analysis software, such as SciPy.

14.3.2 Smoothed ReLU function

Example 14.2 (Smoothed ReLU). The ReLU is defined by

$$r(u) := \begin{cases} u & \text{if } u \geq 0 \\ 0 & \text{otherwise} \end{cases} = u \cdot h(u).$$

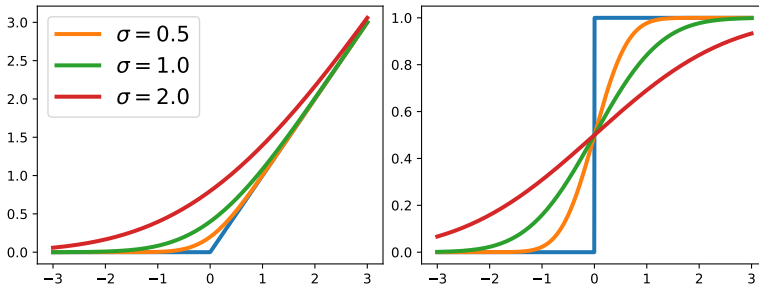


Figure 14.3: Smoothing of the ReLU and Heaviside functions by convolution with a Gaussian kernel, for three values of the width σ .

Similarly to the previous example, we obtain

$$\begin{aligned}
 (r * \kappa_\sigma)(\mu) &= \int_{-\infty}^{\mu} \kappa_\sigma(z) r(\mu - z) dz \\
 &= \int_{-\infty}^{\mu} \kappa_\sigma(z) (\mu - z) dz \\
 &= \mu \int_{-\infty}^{\mu} \kappa_\sigma(z) dz - \int_{-\infty}^{\mu} \kappa_\sigma(z) z dz \\
 &= \mu \Phi_\sigma(\mu) + \sigma^2 \kappa_\sigma(\mu).
 \end{aligned}$$

In the second integral, setting $a := \frac{1}{2\sigma^2}$, we used

$$\int z e^{-az^2} dz = -\frac{1}{2a} \int e^t dt = -\frac{1}{2a} e^t + C = -\frac{1}{2a} e^{-az^2} + C$$

and $t := -az^2 \Rightarrow z dz = -\frac{1}{2a} dt$.

To illustrate differentiation of the convolution, we show how to differentiate the smoothed ReLU.

Example 14.3 (Differentiating the smoothed ReLU). Differentiating the smoothed ReLU from Example 14.2, we obtain

$$(r * \kappa_\sigma)' = (r' * \kappa_\sigma) = h * \kappa_\sigma = \Phi_\sigma.$$

Therefore, unsurprisingly, the derivative of the smoothed ReLU is the smoothed Heaviside step function. Differentiating once again,

we obtain,

$$(r * \kappa_\sigma)'' = (h * \kappa_\sigma)' = (h' * \kappa_\sigma) = \delta * \kappa_\sigma = \kappa_\sigma,$$

where the derivative h' is well-defined almost everywhere. We can arrive at the same result by using that $h * \kappa_\sigma = \Phi_\sigma$ and $\Phi'_\sigma = \kappa_\sigma$, since Φ_σ and κ_σ are the CDF and PDF of the Gaussian with zero mean and σ^2 variance.

14.4 Perturbation of blackbox functions

In this section, we review how to approximately compute a convolution with a kernel and its gradient using Monte-Carlo estimation.

14.4.1 Expectation in a location-scale family

A rather intuitive approach to smooth a function $f : \mathbb{R}^M \rightarrow \mathbb{R}$ is to average its values on an input $\boldsymbol{\mu}$, perturbed by some additive noise $Z \sim p$, for some noise distribution p . This defines the surrogate

$$f_\sigma(\boldsymbol{\mu}) := \mathbb{E}_{Z \sim p}[f(\boldsymbol{\mu} + \sigma Z)].$$

The parameter σ controls the perturbation strength: as $\sigma \rightarrow 0$, we naturally recover f . An equivalent viewpoint is obtained by defining the transformation (change of variables)

$$U := \boldsymbol{\mu} + \sigma Z.$$

We then have

$$U \sim p_{\boldsymbol{\mu}, \sigma},$$

where $p_{\boldsymbol{\mu}, \sigma}$ is the **location-family distribution** generated by the noise distribution p . It is the pushforward distribution of Z through the transformation (see Section 12.4.4). In this notation, the initial noise distribution p is then simply $p = p_{\mathbf{0}, 1}$. The perturbed function can then be expressed from these two perspectives as

$$\begin{aligned} f_\sigma(\boldsymbol{\mu}) &= \mathbb{E}_{Z \sim p_{\mathbf{0}, 1}}[f(\boldsymbol{\mu} + \sigma \cdot Z)] \\ &= \mathbb{E}_{U \sim p_{\boldsymbol{\mu}, \sigma}}[f(U)]. \end{aligned} \tag{14.4}$$

Writing the expectation as the integral of a p.d.f, we naturally recover the smoothing by convolution presented earlier,

$$\begin{aligned} f_\sigma(\boldsymbol{\mu}) &= \int f(\boldsymbol{\mu} + \sigma \mathbf{z}) p_{\mathbf{0},1}(\mathbf{z}) d\mathbf{z} \\ &= f * \kappa_\sigma(\boldsymbol{\mu}), \end{aligned}$$

where we defined the kernel

$$\kappa_\sigma(\mathbf{z}) := p_{\mathbf{0},\sigma}(-\mathbf{z}).$$

In the sequel, we assume that the noise distribution decomposes as

$$p_{\mathbf{0},1}(\mathbf{z}) := \exp(-\nu(\mathbf{z}))/C,$$

where $\nu(\mathbf{z})$ is the **log-density** of the noise distribution and C is a normalization constant. For instance, the Gaussian distribution with diagonal covariance matrix and the corresponding **Gaussian kernel** are obtained with $\nu(\mathbf{z}) = \frac{1}{2}\|\mathbf{z}\|_2^2$ and $C = \sqrt{2\pi}^M$.

Approximation by Monte-Carlo estimation

Instead of approximating the integral above (continuous convolution) with a discrete convolution on a grid, as we did in Section 14.1.3, the expectation perspective suggests that we can estimate $f_\sigma(\boldsymbol{\mu})$ by Monte-Carlo estimation: we simply draw samples from the distribution, evaluate the function at these samples and average. Beyond mere Monte-Carlo estimation, more elaborate approximation schemes are studied in (Chaudhuri and Solar-Lezama, 2010).

14.4.2 Gradient estimation by reparametrization

Provided that the conditions for swapping differentiation and integration hold (see Section 12.1), we have

$$\nabla f_\sigma(\boldsymbol{\mu}) = \mathbb{E}_{\mathbf{Z} \sim p_{\mathbf{0},1}}[\nabla f(\boldsymbol{\mu} + \sigma \cdot \mathbf{Z})]. \quad (14.5)$$

Note that if f is only differentiable almost everywhere, the formula may still hold. For example, if f is the ReLU, then ∇f is the Heaviside step function, and we obtain the correct gradient of f_σ using the formula

above; see Example 14.1. However, if f is not absolutely continuous, the formula may not hold. For example, if f is the Heaviside function, the right-hand side of (14.5) is 0 which does not match the gradient of f_σ ; see again Example 14.1.

From the second expression of f_σ in (14.4), we can see the formula of the gradient in (14.5) as a reparametrization trick $U = \boldsymbol{\mu} + \sigma Z$; see Section 12.4. Namely, we have

$$\begin{aligned}\nabla f_\sigma(\boldsymbol{\mu}) &= \nabla_{\boldsymbol{\mu}} \mathbb{E}_{U \sim p_{\boldsymbol{\mu}, \sigma}}[f(U)] \\ &= \nabla_{\boldsymbol{\mu}} \mathbb{E}_{Z \sim p_{\mathbf{0}, 1}}[f(\boldsymbol{\mu} + \sigma \cdot Z)] \\ &= \mathbb{E}_{Z \sim p_{\mathbf{0}, 1}}[\nabla_{\boldsymbol{\mu}} f(\boldsymbol{\mu} + \sigma \cdot Z)] \\ &= \mathbb{E}_{Z \sim p_{\mathbf{0}, 1}}[\nabla f(\boldsymbol{\mu} + \sigma \cdot Z)].\end{aligned}\tag{14.6}$$

14.4.3 Gradient estimation by SFE, Stein's lemma

In some cases, we may not have access to ∇f or f may not be absolutely continuous and therefore the formula in (14.5) cannot apply. For these cases, we can use the score function estimator (SFE) from Section 12.3. Here, for $f_\sigma(\boldsymbol{\mu}) = \mathbb{E}_{U \sim p_{\boldsymbol{\mu}, \sigma}}[f(U)]$, we obtain

$$\nabla f_\sigma(\boldsymbol{\mu}) = \mathbb{E}_{U \sim p_{\boldsymbol{\mu}, \sigma}}[f(U) \nabla_{\boldsymbol{\mu}} \log p_{\boldsymbol{\mu}, \sigma}(U)].$$

Since the PDF can be written as

$$p_{\boldsymbol{\mu}, \sigma}(\mathbf{u}) = \frac{1}{\sigma} p_{\mathbf{0}, 1}((\mathbf{u} - \boldsymbol{\mu})/\sigma),$$

where

$$p_{\mathbf{0}, 1}(z) := \exp(-\nu(z))/C,$$

we obtain

$$\nabla_{\boldsymbol{\mu}} \log p_{\boldsymbol{\mu}, \sigma}(\mathbf{u}) = \nabla \nu((\mathbf{u} - \boldsymbol{\mu})/\sigma)/\sigma.$$

To summarize, we have shown that

$$\begin{aligned}\nabla f_\sigma(\boldsymbol{\mu}) &= \mathbb{E}_{U \sim p_{\boldsymbol{\mu}, \sigma}}[f(U) \nabla \nu((U - \boldsymbol{\mu})/\sigma)/\sigma] \\ &= \mathbb{E}_{Z \sim p_{\mathbf{0}, 1}}[f(\boldsymbol{\mu} + \sigma \cdot Z) \nabla \nu(Z)/\sigma],\end{aligned}\tag{14.7}$$

where we used the change of variable $Z = (U - \boldsymbol{\mu})/\sigma$. The same technique can also be used if we want to estimate the gradient w.r.t. $\boldsymbol{\theta} = (\boldsymbol{\mu}, \sigma)$ or

if we want to estimate the Jacobian of the expectation of a vector-valued function.

In the particular case of Gaussian noise, since $\nabla \nu(\mathbf{z}) = \mathbf{z}$, we obtain

$$\nabla f_\sigma(\boldsymbol{\mu}) = \mathbb{E}_{Z \sim p_{0,1}}[f(\boldsymbol{\mu} + \sigma \cdot Z)Z/\sigma].$$

This is known as **Stein's lemma**. It should be noted that the above is an unbiased estimator of the gradient of the smoothed function f_σ , but a biased estimator of the gradient of the original function f (assuming that it exists). However, smoothing is usually a good thing, as it can accelerate the convergence of gradient-based algorithms. Computing the gradient of perturbed general programs is studied in detail in (Kreikemeyer and Andelfinger, 2023).

14.4.4 Link between reparametrization and SFE

Using the log-derivative identity, we have for any distribution with differentiable density p

$$\begin{aligned} \mathbb{E}_{Z \sim p}[h(Z)\nabla \log p(Z)] &= \int_{\mathbb{R}^M} h(\mathbf{z}) \left(\frac{\nabla p(\mathbf{z})}{p(\mathbf{z})} \right) p(\mathbf{z}) d\mathbf{z} \\ &= \int_{\mathbb{R}^M} h(\mathbf{z}) \nabla p(\mathbf{z}) d\mathbf{z}. \end{aligned}$$

Using **integration by parts** and assuming that $h(\mathbf{z})p(\mathbf{z})$ goes to zero when $\|\mathbf{z}\| \rightarrow \infty$, we have

$$\int_{\mathbb{R}^M} h(\mathbf{z}) \nabla p(\mathbf{z}) d\mathbf{z} = - \int_{\mathbb{R}^M} p(\mathbf{z}) \nabla h(\mathbf{z}) d\mathbf{z}.$$

We have therefore the identity

$$\mathbb{E}_{Z \sim p}[h(Z)\nabla \log p(Z)] = -\mathbb{E}_{Z \sim p}[\nabla h(Z)].$$

Importantly, contrary to the SFE estimator from Section 12.3, this identity uses gradients with respect to \mathbf{z} , not with respect to the parameters of the distribution. Nevertheless, using the reparametrization

$h(\mathbf{z}) := f(\boldsymbol{\mu} + \sigma \cdot \mathbf{z})$, we have $\nabla h(\mathbf{z}) = \nabla f(\boldsymbol{\mu} + \sigma \cdot \mathbf{z}) \cdot \sigma$ so that

$$\begin{aligned}
 \nabla f_\sigma(\boldsymbol{\mu}) &= \nabla_{\boldsymbol{\mu}} \mathbb{E}_{U \sim p_{\boldsymbol{\mu}, \sigma}} [f(U)] \\
 &= \mathbb{E}_{Z \sim p_{\mathbf{0}, 1}} [\nabla f(\boldsymbol{\mu} + \sigma \cdot Z)] \quad (\text{reparametrization trick}) \\
 &= \mathbb{E}_{Z \sim p_{\mathbf{0}, 1}} [\nabla h(Z) / \sigma] \\
 &= -\mathbb{E}_{Z \sim p_{\mathbf{0}, 1}} [h(Z) \nabla \log p(Z) / \sigma] \\
 &= \mathbb{E}_{Z \sim p_{\mathbf{0}, 1}} [h(Z) \nabla \nu(Z) / \sigma] \\
 &= \mathbb{E}_{Z \sim p_{\mathbf{0}, 1}} [f(\boldsymbol{\mu} + \sigma \cdot Z) \nabla \nu(Z) / \sigma] \quad (\text{score function estimator})
 \end{aligned}$$

Essentially, integration by parts allowed us to convert the reparametrization trick estimator into the SFE estimator. For more applications of integration by parts in machine learning, see Francis Bach’s excellent [blog post](#).

14.4.5 Variance reduction and evolution strategies

As discussed in Chapter 12, the SFE suffers from high variance. We now apply variance reduction techniques to it. To do so, we assume that $\nabla \nu(Z)$ has zero mean for $Z \sim p_{\mathbf{0}, 1}$. This assumption for example holds for Gaussian noise. This assumption implies that

$$\mathbb{E}_{Z \sim p_{\mathbf{0}, 1}} [f(\boldsymbol{\mu}) \nabla \nu(Z) / \sigma] = f(\boldsymbol{\mu}) \mathbb{E}_{Z \sim p_{\mathbf{0}, 1}} [\nabla \nu(Z) / \sigma] = \mathbf{0}$$

and therefore

$$\nabla f_\sigma(\boldsymbol{\mu}) = \mathbb{E}_{Z \sim p_{\mathbf{0}, 1}} [(f(\boldsymbol{\mu} + \sigma \cdot Z) - f(\boldsymbol{\mu})) \nabla \nu(Z) / \sigma]. \quad (14.8)$$

This is an example of **control variate** discussed in Section 12.3. This can be interpreted as using a **finite difference** for computing a directional derivative in the **random direction** Z (see “limit case” below). Inspired by a **central finite difference**, we can also use

$$\nabla f_\sigma(\boldsymbol{\mu}) = \mathbb{E}_{Z \sim p_{\mathbf{0}, 1}} [(f(\boldsymbol{\mu} + \sigma \cdot Z) - f(\boldsymbol{\mu} - \sigma \cdot Z)) \nabla \nu(Z) / (2\sigma)]. \quad (14.9)$$

These estimators have been used as part of blackbox (zero-order) optimization algorithms, such as **evolution strategies** (Salimans *et al.*, 2017) or **random gradient-free optimization** (Nesterov and Spokoiny, 2017). For quadratic functions, it is easy to show that the second estimator achieves lower variance (Recht and Frostig, 2017). The

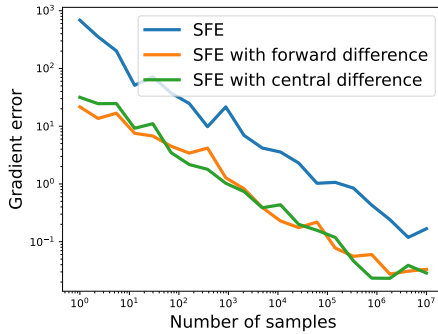


Figure 14.4: Comparison of the score function estimator (SFE) with or without variance reduction for blackbox gradient estimation. We show the error $|\nabla f(\mu) - \nabla f_\sigma(\mu)|$ for $f(u) := u^3$ and $f_\sigma(\mu) := \mathbb{E}[f(\mu + \sigma Z)]$, where $Z \sim \text{Normal}(0, 1)$ and $\sigma := 0.1$. To estimate $\nabla f_\sigma(\mu)$, we compare three estimators: the vanilla SFE Eq. (14.7), the SFE estimator with forward difference (variance reduced) Eq. (14.8), and the SFE estimator with central difference (variance reduced) Eq. (14.9). In all three cases, we approximate the expectation by Monte-Carlo estimation using some number of samples. The variance-reduced estimators not only achieve smaller error, they are also more numerically stable as σ gets smaller.

idea of sampling both Z and $-Z$ simultaneously is called antithetic (Geweke, 1988) or mirrored sampling (Brockhoff *et al.*, 2010). Evolution strategies have also been used to obtain unbiased gradient estimators of partially unrolled computational graphs (Vicol *et al.*, 2021). We empirically compare the SFE with or without variance reduction for blackbox gradient estimation in Fig. 14.4.

14.4.6 Zero-temperature limit

We now discuss the limit case $\sigma \rightarrow 0$. That is, we assume that we do **not** want to perform smoothing and that ∇f exists. We recall that the directional derivative of f at μ in the direction z is

$$\begin{aligned} \partial f(\mu)[z] &= \langle \nabla f(\mu), z \rangle \\ &= \lim_{\sigma \rightarrow 0} [f(\mu + \sigma \cdot z) - f(\mu)] / \sigma. \end{aligned}$$

When $\sigma \rightarrow 0$ and Z follows the standard Gaussian distribution, meaning that $\nabla \nu(z) = z$, Eq. (14.8) therefore becomes

$$\begin{aligned} \nabla f_\sigma(\boldsymbol{\mu}) &= \mathbb{E}_{Z \sim p_{0,1}} [\partial f(\boldsymbol{\mu})[Z] \nabla \nu(Z)] \\ &= \mathbb{E}_{Z \sim p_{0,1}} [\partial f(\boldsymbol{\mu})[Z] Z] \\ &= \mathbb{E}_{Z \sim p_{0,1}} [\langle \nabla f(\boldsymbol{\mu}), Z \rangle Z] \\ &= \mathbb{E}_{Z \sim p_{0,1}} [\nabla f(\boldsymbol{\mu}) Z Z^\top] \\ &= \nabla f(\boldsymbol{\mu}). \end{aligned}$$

This should not be too surprising, as we already know from the convolution perspective that $f_\sigma(\boldsymbol{\mu}) = (f * \kappa_\sigma)(\boldsymbol{\mu}) \rightarrow f(\boldsymbol{\mu})$ when $\sigma \rightarrow 0$. This recovers the randomized forward-mode estimator already presented in Section 8.7.

14.5 Gumbel tricks

14.5.1 The Gumbel distribution

The Gumbel distribution is a distribution frequently used in extreme value theory. As illustrated in Fig. 14.5, we consider the shifted standard Gumbel distribution, whose PDF is defined by

$$p(z) := \exp(-\nu(z)),$$

where

$$\nu(z) := z + \gamma + \exp(-(z + \gamma)),$$

and where $\gamma \approx 0.577$ is Euler's constant. Note that in some formulations, the distribution is not shifted, i.e., γ is not added. If Z is distributed according to the shifted standard Gumbel distribution, we write $Z \sim \text{Gumbel}(0, 1)$.

To obtain a multivariate extension with location-scale parameters $\boldsymbol{\mu}$ and σ , we take M independent random variables $Z := (Z_1, \dots, Z_m)$ and apply the location-scale transform (Section 12.4.1). That is,

$$U \sim \text{Gumbel}(\boldsymbol{\mu}, \sigma) \iff U = \boldsymbol{\mu} + \sigma Z, \quad Z_i \sim \text{Gumbel}(0, 1).$$

As we used shifted standard Gumbel distributions, we naturally get that $\mathbb{E}[U] = \boldsymbol{\mu}$ and $\text{Var}(U) = \sigma^2$. We can use Gumbel noise as an alternative

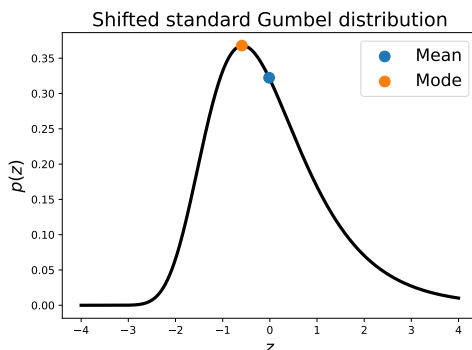


Figure 14.5: We use a shifted definition of the standard Gumbel distribution so that the mean is achieved at $z = 0$, and the mode at $z = -\gamma$. With an unshifted definition, the mean and the mode would be achieved at γ and 0, respectively.

to the Gaussian noise used in Section 14.4. Thankfully, in particular cases, we can compute closed-form expressions of the expectation of perturbed functions.

Remark 14.2 (Link between Gumbel and exponential distribution).

A random variable Z is distributed as $\text{Gumbel}(\mu, 1)$ if and only if $\exp(-Z)$ is distributed as an exponential distribution $\text{Exp}(\exp(\mu - \gamma))$. To see this, one can simply compute the CDF of $\exp(-Z)$ and recognize the CDF of $\text{Exp}(\exp(\mu - \gamma))$. Therefore, when comparing Gumbel distributions, we can use standard properties of the exponential distribution.

Remark 14.3 (Sampling Gumbel noise). If $U \sim \text{Uniform}(0, 1)$, then

$Z := -\log(-\log(U)) - \gamma$ satisfies $Z \sim \text{Gumbel}(0, 1)$, where we recall that we use $\text{Gumbel}(0, 1)$ to denote the shifted standard Gumbel distribution. To see this, note that $\mathbb{P}(Z \leq t) = \mathbb{P}(U \leq \exp(\exp(-(\gamma + t)))) = \exp(\exp(-(\gamma + t)))$, where the last expression matches the CDF of $\text{Gumbel}(0, 1)$.

14.5.2 Perturbed comparison

To start with, the Gumbel distribution can be used to smooth a binary comparison like the greater than or equal operators. Recall that the latter is defined for any $\mu_1, \mu_2 \in \mathbb{R}$ as

$$\text{gt}(\mu_1, \mu_2) := \begin{cases} 1 & \text{if } \mu_1 \geq \mu_2 \\ 0 & \text{if } \mu_1 < \mu_2 \end{cases} = \text{step}(\mu_1 - \mu_2),$$

where step is the Heaviside function. As shown below, by perturbing each variable with Gumbel noise, we recover $\text{logistic}(a - b) = 1/(1 + e^{-(a-b)})$ as an approximation of $\text{step}(a - b)$.

Proposition 14.3 (Gumbel trick for binary variables). Let $Z_1, Z_2 \sim \text{Gumbel}(0, 1)$ be two independent random variables. The difference of their location-scale transform (Section 12.4.1) is distributed according to a logistic distribution (Remark 3.1), i.e.,

$$\mu_1 + \sigma Z_1 - (\mu_2 + \sigma Z_2) \sim \text{Logistic}(\mu_1 - \mu_2, \sigma),$$

for $\mu_1, \mu_2 \in \mathbb{R}$ and $\sigma > 0$. In particular, we have

$$\mathbb{E}_{Z_1, Z_2}[\text{gt}(\mu_1 + \sigma Z_1, \mu_2 + \sigma Z_2)] = \frac{1}{1 + e^{-(\mu_1 - \mu_2)/\sigma}}.$$

Proof. We first derive the CDF of $\mu_1 + \sigma Z_1 - (\mu_2 + \sigma Z_2)$ as

$$\begin{aligned} \mathbb{P}(\mu_1 + \sigma Z_1 - (\mu_2 + \sigma Z_2) \leq t) &= \mathbb{P}(\mu_1/\sigma + Z_1 \leq (\mu_2 + t)/\sigma + Z_2) \\ &= \mathbb{P}\left(e^{-(\mu_1/\sigma + Z_1)} \geq e^{-((\mu_2 + t)/\sigma + Z_2)}\right). \end{aligned}$$

By Remark 14.2, $e^{-(\mu_1/\sigma + Z_1)} \sim \text{Exp}(\exp(\mu_1/\sigma - \gamma))$, and similarly for $e^{-((\mu_2 + t)/\sigma + Z_2)}$. Now one easily shows that if $U \sim \text{Exp}(u)$, $V \sim \text{Exp}(v)$ independent, then $\mathbb{P}(U \leq V) = u/(u + v)$. Hence, we get

$$\begin{aligned} \mathbb{P}(\mu_1 + \sigma Z_1 - (\mu_2 + \sigma Z_2) \leq t) &= \frac{e^{(\mu_2 + t)/\sigma - \gamma}}{e^{(\mu_2 + t)/\sigma - \gamma} + e^{\mu_1/\sigma - \gamma}} \\ &= \frac{1}{1 + e^{-(t - (\mu_1 - \mu_2))/\sigma}}. \end{aligned}$$

We recognize the CDF of the logistic distribution with mean $\mu_1 - \mu_2$ and scale σ , denoted $\text{Logistic}(\mu_1 - \mu_2, \sigma)$. For the last claim, we simply

have that

$$\begin{aligned}\mathbb{E}[\text{gt}(\mu_1 + \sigma Z_1, \mu_2 + \sigma Z_2)] &= \mathbb{E}[\text{step}(\mu_1 + \sigma Z_1 - (\mu_2 + \sigma Z_2))] \\ &= \mathbb{P}(\mu_1 + \sigma Z_1 - (\mu_2 + \sigma Z_2) \geq 0) \\ &= \frac{1}{1 + e^{-(\mu_1 - \mu_2)/\sigma}}.\end{aligned}$$

□

14.5.3 Perturbed argmax

Suppose we want to smooth

$$\mathbf{y}(\mathbf{u}) := \arg \max_{\mathbf{y} \in \{\mathbf{e}_1, \dots, \mathbf{e}_M\}} \langle \mathbf{y}, \mathbf{u} \rangle = \phi(i(\mathbf{u})),$$

where

$$\begin{aligned}i(\mathbf{u}) &:= \arg \max_{i \in [M]} u_i \\ \phi(i) &:= \mathbf{e}_i\end{aligned}$$

with $\phi(i)$ is the one-hot encoding of $i \in [M]$. It turns out that the function $\mathbf{y}(\mathbf{u})$ perturbed using Gumbel noise enjoys a closed form expectation, which is nothing else than the softargmax.

Proposition 14.4 (Gumbel trick for categorical variables). Let us define M independent random variables $Z \sim \text{Gumbel}(\mathbf{0}, 1) \in \mathbb{R}^M$. Then, for $\boldsymbol{\mu} \in \mathbb{R}^M$ and $\sigma > 0$,

$$Y := i(\boldsymbol{\mu} + \sigma \cdot Z) \iff Y \sim \text{Categorical}(\text{softargmax}(\boldsymbol{\mu}/\sigma)).$$

Moreover, we have

$$\begin{aligned}\mathbf{y}_\sigma(\boldsymbol{\mu}) &= \mathbb{E}_Z[\mathbf{y}(\boldsymbol{\mu} + \sigma \cdot Z)] \\ &= \mathbb{E}_Y[\phi(Y)] \\ &= \text{softargmax}(\boldsymbol{\mu}/\sigma).\end{aligned}$$

Proof. For $k \in [M]$, we have that

$$\begin{aligned}\mathbb{P}(Y = k) &= \mathbb{P}\left(\arg \max_{i \in [M]} \{\mu_i + \sigma Z_i\} = k\right) \\ &= \mathbb{P}\left(\arg \min_{i \in [M]} \{e^{-\mu_i/\sigma - Z_i}\} = k\right)\end{aligned}$$

By Remark 14.2, we have that $e^{-\mu_i/\sigma - Z_i} \sim \text{Exp}(\exp(\mu_i/\sigma - \gamma))$. One easily verifies as an exercise, that, for U_1, \dots, U_M independent exponential variables with parameters u_1, \dots, u_m , we have $\mathbb{P}(\arg \min_{i \in [M]} \{U_i\} = k) = u_k / \sum_{i=1}^M u_i$. Hence, we get

$$\mathbb{P}(Y = k) = \frac{\exp(\mu_k/\sigma)}{\sum_{i=1}^M \exp(\mu_i/\sigma)},$$

that is,

$$Y \sim \text{Categorical}(\text{softargmax}(\boldsymbol{\mu}/\sigma)).$$

The last claim follows from the distribution of Y and the definition of ϕ . \square

14.5.4 Perturbed max

A similar result holds if we now wish to perturb the max instead of the argmax.

Proposition 14.5 (Link to log-sum-exp). Let us define M independent random variables $Z \sim \text{Gumbel}(\mathbf{0}, 1) \in \mathbb{R}^M$ and

$$f(\mathbf{u}) := \max_{i \in [M]} u_i.$$

Then, for $\boldsymbol{\mu} \in \mathbb{R}^M$ and $\sigma > 0$,

$$V := f(\boldsymbol{\mu} + \sigma \cdot Z) \iff V \sim \text{Gumbel}(\sigma \text{LSE}(\boldsymbol{\mu}/\sigma), \sigma).$$

Moreover, we have

$$f_\sigma(\boldsymbol{\mu}) := \mathbb{E}_Z[f(\boldsymbol{\mu} + \sigma \cdot Z)] = \mathbb{E}_V[V] = \sigma \cdot \text{LSE}(\boldsymbol{\mu}/\sigma).$$

Proof. We derive the CDF of $f(\boldsymbol{\mu} + \sigma \cdot Z)$ as

$$\mathbb{P}(\max_{i \in [M]} \{\mu_i + \sigma Z_i\} \leq t) = \mathbb{P}\left(\min_{i \in [M]} \{e^{-(\mu_i/\sigma - Z_i)}\} \geq e^{-t/\sigma}\right)$$

We have $e^{-(\mu_i/\sigma - Z_i)} \sim \text{Exp}(\exp(\mu_i/\sigma_i) - \gamma)$ and for U_1, \dots, U_M independent exponential random variables with parameters u_i , we have $\min_{i \in [M]} U_i \sim \text{Exp}(\sum_{i=1}^M u_i)$. Hence,

$$\begin{aligned} \mathbb{P}\left(\max_{i \in [M]} \{\mu_i + \sigma Z_i\} \leq t\right) &= \exp\left(-\sum_{i=1}^M (\exp(\mu_i/\sigma - \gamma)) \exp(-t/\sigma)\right) \\ &= \exp(-\exp(-(t - \sigma \text{LSE}(\boldsymbol{\mu}/\sigma))/\sigma - \gamma)). \end{aligned}$$

We recognize the CDF of the shifted Gumbel distribution with location-scale parameters $\sigma \text{LSE}(\boldsymbol{\mu}/\sigma)$ and σ . \square

For further reading on the Gumbel trick, see Tim Vieira's great [blog](#).

14.5.5 Gumbel trick for sampling

The Gumbel trick is also useful in its own right for **sampling** without computing the normalization constant of the softargmax. Indeed, Proposition 14.4 ensures that if Z is Gumbel noise, then Y is distributed according to $\text{Categorical}(\text{softargmax}(\boldsymbol{\mu}/\sigma))$. Computing the arg-maximum, as required to compute Y , can be done in one pass. Therefore, we obtain a one-pass algorithm to sample directly from the logits $\boldsymbol{\mu}$, without explicitly computing the probabilities $\text{softargmax}(\boldsymbol{\mu}/\sigma)$.

One may wonder whether such trick could also be used with the normal distribution. Unfortunately, there is no closed form in this case because it would require integrating the CDF of the maximum of $M - 1$ Gaussian distributions. However, other tricks can be defined such as using Weibull distributions, see Balog *et al.* (2017).

14.5.6 Perturb-and-MAP

Previously, we discussed the Gumbel trick in the classification setting, where $\mathcal{Y} = [M]$. In the structured prediction setting, outputs are

typically embedded in \mathbb{R}^M but the output space is very large. That is, $\mathcal{Y} \subseteq \mathbb{R}^M$ but $|\mathcal{Y}| \gg M$. Structured outputs are then decoded using a maximum a-posteriori (MAP) oracle

$$\begin{aligned} f(\mathbf{u}) &:= \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, \mathbf{u} \rangle \\ \mathbf{y}(\mathbf{u}) &:= \arg \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, \mathbf{u} \rangle. \end{aligned}$$

For this setting, the perturbed versions of f and \mathbf{y} ,

$$\begin{aligned} f_\sigma(\boldsymbol{\mu}) &:= \mathbb{E}_Z[f(\boldsymbol{\mu} + \sigma \cdot Z)] \\ \mathbf{y}_\sigma(\boldsymbol{\mu}) &:= \mathbb{E}_Z[\mathbf{y}(\boldsymbol{\mu} + \sigma \cdot Z)], \end{aligned}$$

no longer enjoy a closed form in general. However, we can approximate them using Monte-carlo estimation. For the gradient of $\nabla f_\sigma(\boldsymbol{\mu})$, two estimators exist (Abernethy *et al.*, 2016; Berthet *et al.*, 2020).

Proposition 14.6 (Gradient of perturbed max). Let $\mathcal{Y} \subseteq \mathbb{R}^M$ and Z be noise with density

$$p_{0,1}(\mathbf{z}) := \exp(-\nu(\mathbf{z}))/C.$$

Then, $f_\sigma(\boldsymbol{\mu})$ is smooth, and its gradient is given by

$$\begin{aligned} \nabla f_\sigma(\boldsymbol{\mu}) &= \mathbb{E}_Z[\mathbf{y}(\boldsymbol{\mu} + \sigma \cdot Z)] \\ &= \mathbb{E}_Z[f(\boldsymbol{\mu} + \sigma \cdot Z) \nabla \nu(Z) \sigma] \\ &\in \text{conv}(\mathcal{Y}). \end{aligned}$$

We therefore have $\nabla f_\sigma(\boldsymbol{\mu}) = \mathbf{y}_\sigma(\boldsymbol{\mu})$.

The first estimator is simply a consequence of the reparametrization trick seen in Eq. (14.6) and of $\mathbf{y} = \nabla f$, which follows from Danskin's theorem (see Section 11.2). The second estimator is just SFE seen in Eq. (14.7). The first estimator usually has lower variance, as it uses more information, namely that $\mathbf{y} = \nabla f$.

The Jacobian of $\mathbf{y}_\sigma(\boldsymbol{\mu})$ also has two estimators (Abernethy *et al.*, 2016; Berthet *et al.*, 2020).

Proposition 14.7 (Jacobian of perturbed argmax). Under the same notation as in Proposition 14.6, we have

$$\begin{aligned}\partial \mathbf{y}_\sigma(\boldsymbol{\mu}) &= \mathbb{E}_Z \left[\mathbf{y}(\boldsymbol{\mu} + \sigma Z) \nabla \nu(Z)^\top / \sigma \right] \\ &= \mathbb{E}_Z \left[f(\boldsymbol{\mu} + \sigma Z) \left(\nabla \nu(Z) \nabla \nu(Z)^\top - \nabla^2 \nu(Z) \right) / \sigma^2 \right].\end{aligned}$$

The first estimator uses SFE. The second estimator is obtained by differentiating through

$$\mathbf{y}_\sigma(\boldsymbol{\mu}) = \nabla f_\sigma(\boldsymbol{\mu}) = \mathbb{E}_Z [f(\boldsymbol{\mu} + \sigma \cdot Z) \nabla \nu(Z) / \sigma].$$

The first estimator usually has lower variance. Note that we cannot use the reparametrization trick this time, since \mathbf{y} is discontinuous, contrary to f .

Link between perturbation and regularization

As shown in (Berthet *et al.*, 2020, Proposition 2.2), assuming \mathcal{Y} is a convex polytope with non-empty interior and p has a strictly positive density, the function

$$f_\sigma(\boldsymbol{\mu}) := \mathbb{E}_Z [f(\boldsymbol{\mu} + \sigma \cdot Z)] = \mathbb{E}_Z [\max_{\mathbf{y} \in \mathcal{Y}} \langle \boldsymbol{\mu} + \sigma \cdot Z, \mathbf{y} \rangle]$$

is strictly convex and its convex conjugate $f_\sigma^*(\mathbf{y})$ is Legendre-type. We can therefore rewrite $f_\sigma(\boldsymbol{\mu})$ from the regularization perspective as

$$f_\sigma(\boldsymbol{\mu}) = \max_{\mathbf{y} \in \mathcal{Y}} \langle \boldsymbol{\mu}, \mathbf{y} \rangle - f_\sigma^*(\mathbf{y}).$$

and $\nabla f_\sigma(\boldsymbol{\mu}) = \mathbf{y}_\sigma(\boldsymbol{\mu})$ is a **mirror map**, a one-to-one mapping from \mathbb{R}^M to the interior of \mathcal{Y} . Unfortunately, $f_\sigma^*(\mathbf{y})$ does not enjoy a closed form in general. Conversely, does any regularization has a corresponding noise distribution? The reciprocal is not true.

14.5.7 Gumbel-softargmax

Suppose we want to smooth out the **composition** $h(\mathbf{u}) := g(\mathbf{y}(\mathbf{u}))$ between some function $g: \{\mathbf{e}_1, \dots, \mathbf{e}_M\} \rightarrow \mathbb{R}$ and the argmax

$$\mathbf{y}(\mathbf{u}) := \arg \max_{\mathbf{y} \in \{\mathbf{e}_1, \dots, \mathbf{e}_M\}} \langle \mathbf{y}, \mathbf{u} \rangle.$$

We can do so by

$$h_\sigma(\boldsymbol{\mu}) := \mathbb{E}_Z [g(\mathbf{y}(\boldsymbol{\mu} + \sigma Z))].$$

This is useful for instance to compute the expectation of a loss (instead of the loss of an expectation). To compute the gradient of $h_\sigma(\boldsymbol{\mu})$, we can readily use the SFE. However, we saw that it suffers from high variance. Unfortunately, we cannot swap differentiation and integration (expectation) here, since $\mathbf{y}(\mathbf{u})$ is a discontinuous function. See Section 12.1 for more details regarding differentiation under the integral sign.

The key idea of the Gumbel-softargmax (Jang *et al.*, 2016; Maddison *et al.*, 2016) is to replace $\mathbf{y}(\mathbf{u})$ with a softargmax (with temperature parameter τ) to define

$$h_{\sigma,\tau}(\boldsymbol{\mu}) := \mathbb{E}_Z [g(\text{softargmax}_\tau(\boldsymbol{\mu} + \sigma Z))].$$

Since the softargmax is a regularized argmax, we can see the Gumbel-softargmax approach as using **both** regularization and perturbation. Note that the approach is also known as Gumbel-softmax. We use the name Gumbel-softargmax for consistency with the terminology of this book.

The key benefit is that we can now swap differentiation and integration (expectation) to get an unbiased estimator of $\nabla h_{\sigma,\tau}(\boldsymbol{\mu})$. However, this will be a **biased** estimator of $\nabla h_\sigma(\boldsymbol{\mu})$, the amount of bias being controlled by the temperature τ . In particular, in the limit case $\tau \rightarrow 0$, we have $h_{\sigma,\tau}(\boldsymbol{\mu}) \rightarrow h_\sigma(\boldsymbol{\mu})$. One caveat, however, is that the function g needs to be well defined on Δ^M , instead of $\{\mathbf{e}_1, \dots, \mathbf{e}_M\}$.

The use of the softargmax transformation defines a continuous distribution (Jang *et al.*, 2016; Maddison *et al.*, 2016), that we now explain with $\sigma = 1$.

Proposition 14.8 (Gumbel-softargmax / Concrete distributions). Let us define the **continuous** random variable

$$T := \text{softargmax}_\tau(\boldsymbol{\mu} + Z) \in \Delta^M,$$

where Z is a Gumbel random variable. Then T is distributed

according to a distribution with density

$$p_{\boldsymbol{\mu}, \tau}(\mathbf{t}) := \Gamma(M) \tau^{M-1} \left(\sum_{i=1}^M \frac{\pi_i}{t_i^\tau} \right)^{-M} \prod_{i=1}^M \frac{\pi_i}{t_i^{\tau+1}},$$

where $\boldsymbol{\pi} := \text{softargmax}(\boldsymbol{\mu})$.

We can extend the Gumbel softargmax to the structured setting by replacing

$$\mathbf{y}(\mathbf{u}) := \arg \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, \mathbf{u} \rangle,$$

with its regularized variant (Paulus *et al.*, 2020). Similarly as before, one caveat is that g needs to be well defined on $\text{conv}(\mathcal{Y})$ instead of \mathcal{Y} . Moreover, regularizing \mathbf{y} is not always easy computationally.

14.6 Summary

- We studied smoothing techniques based on function convolution with a kernel. Due to the commutativity of the convolution, we can alternatively see these as the expectation of the function, perturbed with noise, assuming the kernel corresponds to the PDF of some noise distribution.
- Their gradients can be estimated using the path gradient estimator (PGE) or score function estimator (SFE), depending on whether the gradient of the original function is available or not.
- We saw that Stein’s lemma is a special case of SFE used with Gaussian noise. The so-called “evolution strategies” are just a variant of that with variance reduction and can be interpreted as randomized finite difference.
- When using Gumbel noise, we were able to derive closed-form expressions for the expectation in specific cases: perturbed comparison, perturbed argmax and perturbed max.
- We also studied the connections between smoothing by optimization and smoothing by integration. Infimal convolution is the counterpart of convolution, and the Legendre-Fenchel transform

is the counterpart of Fourier and Laplace's transforms. Infimal convolution uses a min-plus algebra in the log domain, while the convolution uses a sum-product algebra in the exponential domain.

Part V

Optimizing differentiable programs

15

Optimization basics

15.1 Objective functions

Consider a function L , for example evaluating the error or “loss” $L(\mathbf{w})$ achieved by a model with parameters $\mathbf{w} \in \mathcal{W}$, where $\mathcal{W} = \mathbb{R}^P$. To find the best possible model parameterization, we seek to minimize $L(\mathbf{w})$, that is, to compute approximately

$$L^* := \inf_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}),$$

assuming that the infimum exists (i.e., $L(\mathbf{w})$ is lower bounded). We will denote a solution, if it exists, by

$$\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}) := \left\{ \mathbf{w} \in \mathcal{W} : L(\mathbf{w}) = \min_{\mathbf{w}' \in \mathcal{W}} L(\mathbf{w}') \right\}.$$

In general, an analytical solution is not available and computing such a minimum approximately requires an optimization algorithm. An optimization algorithm is an iterative procedure, which, starting from an initial point \mathbf{w}^0 , outputs after t iterations a point \mathbf{w}^t that approximates the minimum of L up to some accuracy ε , i.e.,

$$L(\mathbf{w}^t) - L^* \leq \varepsilon. \tag{15.1}$$

15.2 Oracles

To produce iterates $\mathbf{w}^1, \mathbf{w}^2, \dots$ that converge to a minimum, the algorithm naturally needs to have access to information about L . For example, the algorithm needs a priori to be able to evaluate L to know if it decreased its value or not. Such information about the function is formalized by the notion of **oracles** (Nemirovski and Yudin, 1983). Formally, oracles are procedures that an algorithm can call to access information about the objective $L(\mathbf{w})$ at any given point $\mathbf{w} \in \mathcal{W}$. We usually mainly consider the following three oracles.

- **Zero-order oracle:** evaluating the function $L(\mathbf{w}) \in \mathbb{R}$.
- **First-order oracle:** evaluating the gradient $\nabla L(\mathbf{w}) \in \mathcal{W}$ for L differentiable.
- **Second-order oracle:** evaluating the Hessian matrix $\nabla^2 L(\mathbf{w})$, or evaluating the Hessian-vector product (HVP) $\nabla^2 L(\mathbf{w})\mathbf{v} \in \mathcal{W}$, for L twice differentiable and any vector $\mathbf{v} \in \mathcal{W}$.

Given an oracle \mathcal{O} for a function L , we can formally define an optimization algorithm as a procedure which computes the next iterate as a function of all past and current information. Formally, an algorithm \mathcal{A} builds a sequence $\mathbf{w}^1, \dots, \mathbf{w}^t$ from a starting point \mathbf{w}^0 as

$$\mathbf{w}^{t+1} := \mathcal{A}(\mathbf{w}^0, \dots, \mathbf{w}^t, \mathcal{O}(\mathbf{w}^0), \dots, \mathcal{O}(\mathbf{w}^t), \boldsymbol{\lambda}),$$

where $\boldsymbol{\lambda} \in \Lambda \subseteq \mathbb{R}^Q$ encapsulates some hyperparameters of the algorithm, such as the stepsize. Oftentimes, algorithms build the next iterate simply from the information collected at the current iterate, without using all past iterates. That is, they take the form $\mathbf{w}^{t+1} = \mathcal{A}(\mathbf{w}^t, \mathcal{O}(\mathbf{w}^t), \boldsymbol{\lambda})$. A classical example is the gradient descent algorithm, that uses a first-order oracle to compute iterates of the form

$$\mathbf{w}^{t+1} := \mathbf{w}^t - \gamma \nabla L(\mathbf{w}^t),$$

where the stepsize γ is a hyperparameter of the algorithm. The notion of oracle therefore delineates different classes of algorithms. For instance, we may consider zero-order algorithms or first-order algorithms.

15.3 Variational perspective of optimization algorithms

One of the most basic optimization algorithms is the **proximal point** method, which produces \mathbf{w}^{t+1} from \mathbf{w}^t by

$$\mathbf{w}^{t+1} := \arg \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}) + \frac{1}{2\gamma} \|\mathbf{w} - \mathbf{w}^t\|_2^2.$$

In words, the next iterate is produced by solving a trade-off between minimizing the function L and staying close to \mathbf{w}^t . Unfortunately, the optimization problem involved in performing this parameter update is as difficult as the original optimization problem, making the proximal point method impractical.

As we shall see in Chapter 16 and Chapter 17, many optimization algorithms can be seen as an approximation of the proximal point method, in the sense that they solve

$$\mathbf{w}^{t+1} := \arg \min_{\mathbf{w} \in \mathcal{W}} \tilde{L}(\mathbf{w}, \mathbf{w}^t) + \frac{1}{2\gamma} \|\mathbf{w} - \mathbf{w}^t\|_2^2.$$

or more generally

$$\mathbf{w}^{t+1} := \arg \min_{\mathbf{w} \in \mathcal{W}} \tilde{L}(\mathbf{w}, \mathbf{w}^t) + \frac{1}{\gamma} d(\mathbf{w}, \mathbf{w}^t),$$

where $\tilde{L}(\mathbf{w}, \mathbf{w}^t)$ is an approximation of $L(\mathbf{w})$ around \mathbf{w}^t and $d(\mathbf{w}, \mathbf{w}^t)$ is some form of distance between \mathbf{w} and \mathbf{w}^t . Different choices of \tilde{L} and d lead to different optimization algorithms, and to different trade-offs.

15.4 Classes of functions

When studying algorithms theoretically, stronger results can often be stated by restricting to certain classes of functions. We already covered continuous and differentiable functions in Chapter 2. We review a few important other classes in this section.

15.4.1 Lipschitz functions

Lipschitz continuity is a stronger form of continuity. Intuitively, a Lipschitz continuous function is limited in how fast it can change.

Definition 15.1 (Lipschitz-continuous functions). A function $g: \mathcal{W} \rightarrow \mathcal{F}$ is β -Lipschitz continuous if for all $\mathbf{w}, \mathbf{v} \in \mathcal{W}$

$$\|g(\mathbf{w}) - g(\mathbf{v})\|_2 \leq \beta \|\mathbf{w} - \mathbf{v}\|_2.$$

Note that the definition is valid even for vector-valued functions.

With respect to arbitrary norms

Thanks to dual norms reviewed in Section 18.1, we can state a more general definition of Lipschitz continuity based on arbitrary norms, instead of the 2-norm. Moreover, we may consider Lipschitz-continuity over a subset of the input domain.

Definition 15.2 (Lipschitz continuous functions w.r.t. a norm). A function $g: \mathcal{W} \rightarrow \mathcal{F}$ is said to be β -Lipschitz w.r.t. a norm $\|\cdot\|$ over a set $\mathcal{C} \subseteq \mathcal{W}$ if for all $\mathbf{w}, \mathbf{v} \in \mathcal{C}$

$$\|g(\mathbf{w}) - g(\mathbf{v})\|_* \leq \beta \|\mathbf{w} - \mathbf{v}\|.$$

When $\|\cdot\| = \|\cdot\|_2$, we recover Definition 15.1, since the 2-norm is dual to itself.

15.4.2 Smooth functions

A differentiable function L is said to be β -smooth if its gradients are β -Lipschitz continuous. Setting $g(\mathbf{w}) = \nabla L(\mathbf{w})$ in Definition 15.1, we obtain the following definition.

Definition 15.3 (Smooth functions). A differentiable function $L: \mathcal{W} \rightarrow \mathbb{R}$ is β -smooth for $\beta > 0$ if for all $\mathbf{w}, \mathbf{v} \in \mathcal{W}$

$$\|\nabla L(\mathbf{w}) - \nabla L(\mathbf{v})\|_2 \leq \beta \|\mathbf{w} - \mathbf{v}\|_2.$$

Smoothness ensures that the information provided by the gradient at some \mathbf{w} is meaningful in a neighborhood of \mathbf{w} , since its variations are upper-bounded. If the variations were not bounded, the gradient at \mathbf{v}

arbitrarily close to \mathbf{w} could drastically change, rendering the information provided by a first-order oracle potentially useless.

Smoothness of a function can be interpreted as having a quadratic upper bound on the function as formalized below.

Proposition 15.1 (Smooth functions). If a differentiable function $L : \mathcal{W} \rightarrow \mathbb{R}$ is β -smooth then for all $\mathbf{w}, \mathbf{v} \in \mathcal{W}$,

$$|L(\mathbf{w}) - L(\mathbf{v}) + \langle \nabla L(\mathbf{v}), \mathbf{w} - \mathbf{v} \rangle| \leq \frac{\beta}{2} \|\mathbf{w} - \mathbf{v}\|_2^2.$$

In particular, we have

$$L(\mathbf{w}) \leq L(\mathbf{v}) + \langle \nabla L(\mathbf{v}), \mathbf{w} - \mathbf{v} \rangle + \frac{\beta}{2} \|\mathbf{w} - \mathbf{v}\|_2^2.$$

Proof. This is shown by bounding $|L(\mathbf{v}) - L(\mathbf{w}) - \langle \nabla L(\mathbf{w}), \mathbf{v} - \mathbf{w} \rangle|$ using the integral representation of the objective along $\mathbf{w} - \mathbf{v}$, i.e., $|L(\mathbf{v}) - L(\mathbf{w}) - \langle \nabla L(\mathbf{w}), \mathbf{v} - \mathbf{w} \rangle| = |\int_0^1 \langle \nabla L(\mathbf{w} + s(\mathbf{v} - \mathbf{w})), \mathbf{v} - \mathbf{w} \rangle ds - \langle \nabla L(\mathbf{w}), \mathbf{v} - \mathbf{w} \rangle| \leq \int_0^1 \|\nabla L(\mathbf{w} + s(\mathbf{v} - \mathbf{w})) - \nabla L(\mathbf{w})\|_2 \|\mathbf{v} - \mathbf{w}\|_2 ds \leq L\|\mathbf{w} - \mathbf{v}\|_2^2/2$, where the last inequality follows from the smoothness assumption and standard integration. \square

In other words, $L(\mathbf{w})$ is upper-bounded and lower-bounded around \mathbf{v} by a quadratic function of \mathbf{w} . We will see in Section 16.1 that this characterization gives rise to a variational perspective on gradient descent.

With respect to arbitrary norms

We can generalize the definition of smoothness in Definition 15.3 to arbitrary norms.

Definition 15.4 (Smooth functions w.r.t. a norm). A function $L : \mathcal{W} \rightarrow \mathbb{R}$ is β -smooth w.r.t. a norm $\|\cdot\|$ over a set \mathcal{C} if for all $\mathbf{w}, \mathbf{v} \in \mathcal{C}$

$$\|\nabla L(\mathbf{w}) - \nabla L(\mathbf{v})\|_* \leq \frac{\beta}{2} \|\mathbf{w} - \mathbf{v}\|.$$

An equivalent characterization, generalizing Proposition 15.1 to arbitrary norms, is given below (see, e.g. Beck (2017, Theorem 5.8)).

Proposition 15.2 (Smooth functions w.r.t. a norm). If a differentiable function $L : \mathcal{W} \rightarrow \mathbb{R}$ is β -smooth w.r.t. a norm $\|\cdot\|$ over a set \mathcal{C} , then for all $\mathbf{w}, \mathbf{v} \in \mathcal{C}$

$$\underbrace{|L(\mathbf{w}) - L(\mathbf{v}) - \langle \nabla L(\mathbf{v}), \mathbf{w} - \mathbf{v} \rangle|}_{B_L(\mathbf{w}, \mathbf{v})} \leq \frac{\beta}{2} \|\mathbf{w} - \mathbf{v}\|^2,$$

where B_f is the Bregman divergence generated by f (Definition 18.2).

15.4.3 Convex functions

A convex function is a function such that its value on the average of two or more points is smaller than the average of the values of the functions at these points. This is illustrated in Figure 15.2 and formalized below.

Definition 15.5 (Convex functions). A function $L : \mathcal{W} \rightarrow \mathbb{R}$ is said to be **convex** if for all $\mathbf{w}, \mathbf{v} \in \mathcal{W}$ and $\tau \in [0, 1]$

$$L(\tau \mathbf{w} + (1 - \tau) \mathbf{v}) \leq \tau L(\mathbf{w}) + (1 - \tau) L(\mathbf{v}).$$

The function L is **strictly convex** if the above inequality is strict for all $\mathbf{w} \neq \mathbf{v}$.

The above characterization can easily be generalized to multiple points. Namely, for $\mathbf{w}_1, \dots, \mathbf{w}_n \in \mathcal{W}$ and $\tau_1, \dots, \tau_n \geq 0$ such that $\sum_{i=1}^n \tau_i = 1$ (that is, τ_1, \dots, τ_n defines a probability distribution over $[n]$), we have if L is convex that

$$L\left(\sum_{i=1}^n \tau_i \mathbf{w}_i\right) \leq \sum_{i=1}^n \tau_i L(\mathbf{w}_i).$$

The point $\sum_{i=1}^n \tau_i \mathbf{w}_i$ is called a convex combination. This can be seen as comparing the function at the average point to the average of the values at these points and can further be generalized to any random variable.

Proposition 15.3 (Jensen's inequality). A function $L: \mathcal{W} \rightarrow \mathbb{R}$ is convex if it satisfies **Jensen's inequality**, that is, for any random variable W on \mathcal{W} ,

$$L(\mathbb{E}[W]) \leq \mathbb{E}[L(W)],$$

provided that the expectations are well-defined.

If the function considered is differentiable, an alternative characterization of convexity is to observe how linear approximations of the function lower bound the function. This is illustrated in Figure 15.2 and formalized below.

Definition 15.6 (Convex differentiable functions). A differentiable function $L: \mathcal{W} \rightarrow \mathbb{R}$ is convex if and only if for all $\mathbf{w}, \mathbf{v} \in \mathcal{W}$

$$L(\mathbf{v}) \geq L(\mathbf{w}) + \langle \nabla L(\mathbf{w}), \mathbf{v} - \mathbf{w} \rangle.$$

The function L is strictly convex if and only if the above inequality is strict for any $\mathbf{w} \neq \mathbf{v}$.

The above characterization pinpoints the relevance of convex function in optimization: if we can find a point $\hat{\mathbf{w}}$ with null gradient, then we know that we have found the minimum as we have

$$\nabla L(\hat{\mathbf{w}}) = \mathbf{0} \implies \forall \mathbf{v} \in \mathbb{R}^P, L(\mathbf{v}) \geq L(\hat{\mathbf{w}}) \implies L(\hat{\mathbf{w}}) = L^*.$$

This means that by having access to the gradient of the function or an approximation thereof, we have access to a sufficient criterion to know whether we found a global minimum. In the case of a gradient descent on a smooth function, convexity ensures convergence to a minimum at a sublinear rate as detailed below.

Finally, if the function is twice differentiable, convexity of a function can be characterized in terms of the Hessian of the function.

Proposition 15.4 (Convex twice differentiable functions). A twice differentiable function $L: \mathcal{W} \rightarrow \mathbb{R}$ is convex if and only if its Hessian is positive semi-definite,

$$\forall \mathbf{w} \in \mathcal{W}, \nabla^2 L(\mathbf{w}) \succeq 0, \text{ i.e., } \forall \mathbf{w}, \mathbf{v} \in \mathcal{W}, \langle \mathbf{v}, \nabla^2 L(\mathbf{w}) \mathbf{v} \rangle \geq 0.$$

The function L is strictly convex if and only if the Hessian is positive-definite, $\forall \mathbf{w} \in \mathcal{W}, \nabla^2 L(\mathbf{w}) \succ 0$, i.e., $\forall \mathbf{w}, \mathbf{v} \in \mathcal{W}, \langle \mathbf{v}, \nabla^2 L(\mathbf{w}) \mathbf{v} \rangle > 0$.

15.4.4 Strongly-convex functions

Convexity can also be strengthened by considering μ -strongly convex functions.

Definition 15.7 (Strongly-convex functions). A function $L : \mathcal{W} \rightarrow \mathbb{R}$ is μ -strongly convex for $\mu > 0$ if for all $\mathbf{w}, \mathbf{v} \in \mathcal{W}$ and $\tau \in [0, 1]$

$$L(\tau \mathbf{w} + (1 - \tau) \mathbf{v}) \leq \tau L(\mathbf{w}) + (1 - \tau) L(\mathbf{v}) - \frac{\mu}{2} \tau (1 - \tau) \|\mathbf{w} - \mathbf{v}\|_2^2.$$

A differentiable function L is μ -strongly convex if and only if for all $\mathbf{w}, \mathbf{v} \in \mathcal{W}$

$$L(\mathbf{v}) \geq L(\mathbf{w}) + \langle \nabla L(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle + \frac{\mu}{2} \|\mathbf{w} - \mathbf{v}\|_2^2.$$

A twice differentiable function is μ -strongly convex if and only if its Hessian satisfies

$$\forall \mathbf{w} \in \mathcal{W}, \nabla^2 L(\mathbf{w}) \succeq \mu \mathbf{I}, \text{ i.e., } \forall \mathbf{w}, \mathbf{v} \in \mathcal{W}, \langle \mathbf{v}, \nabla^2 L(\mathbf{w}) \mathbf{v} \rangle \geq \mu \|\mathbf{v}\|_2^2.$$

The characterization of strong convexity for differentiable functions states that $L(\mathbf{w})$ is lower-bounded by a quadratic. This enables the design of linearly convergent algorithms as explained later. We naturally have the implications

$$L \text{ strongly convex} \implies L \text{ strictly convex} \implies L \text{ convex}.$$

With respect to arbitrary norms

A function can be strongly convex w.r.t. an arbitrary norm, simply by replacing the 2-norm in Definition 15.7 with that norm. For differentiable strongly convex functions, we have the following alternative characterization, generalizing Definition 15.7 to arbitrary norms.

Proposition 15.5 (Differentiable strongly-convex functions). If a differentiable function $L : \mathcal{W} \rightarrow \mathbb{R}$ is μ -strongly convex w.r.t. a norm $\|\cdot\|$ over a set \mathcal{C} , then for all $\mathbf{w}, \mathbf{v} \in \mathcal{C}$

$$\frac{\mu}{2} \|\mathbf{w} - \mathbf{v}\|^2 \leq \underbrace{L(\mathbf{w}) - L(\mathbf{v}) - \langle \nabla L(\mathbf{v}), \mathbf{w} - \mathbf{v} \rangle}_{B_L(\mathbf{w}, \mathbf{v})}.$$

Obviously, if a function L is μ -strongly convex, then, λL is $(\mu\lambda)$ -strongly convex. Because all norms are equivalent, if a function is strongly convex w.r.t. a norm, it is also strongly-convex w.r.t. another norm. However, stating the norm w.r.t. which strong convexity holds can lead to better constant μ (the higher, the better in terms of convergence rates of, e.g., a gradient descent). We also emphasize that it is important to mention over which set strong convexity holds. We give examples below.

Example 15.1 (Strongly convex functions). The function $f(\mathbf{u}) = \frac{1}{2} \|\mathbf{u}\|_2^2$ is 1-strongly convex w.r.t. $\|\cdot\|_2$ over \mathbb{R}^M .

The function $f(\mathbf{u}) = \langle \mathbf{u}, \log \mathbf{u} \rangle$ is 1-strongly convex w.r.t. $\|\cdot\|_1$ over Δ^M . Applying Proposition 15.5, we obtain for all $\mathbf{p}, \mathbf{q} \in \Delta^M$

$$\frac{1}{2} \|\mathbf{p} - \mathbf{q}\|_1^2 \leq B_f(\mathbf{p}, \mathbf{q}) = \text{KL}(\mathbf{p}, \mathbf{q}),$$

which is known as **Pinsker's inequality**. We empirically verify the inequality in Fig. 15.1.

More generally, $f(\mathbf{u})$ is $\frac{1}{\mu}$ -strongly convex w.r.t. $\|\cdot\|_1$ over any bounded set $\mathcal{C} \subset \mathbb{R}_+^M$ such that $\mu = \sup_{\mathbf{u} \in \mathcal{C}} \|\mathbf{u}\|_1$ (Blondel, 2019). However, it is not strongly convex over \mathbb{R}_+^M , as it is not bounded.

15.4.5 Nonconvex functions

In general, the minimum of a function necessarily has a null gradient, that is,

$$\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}) \implies \nabla L(\mathbf{w}^*) = 0.$$

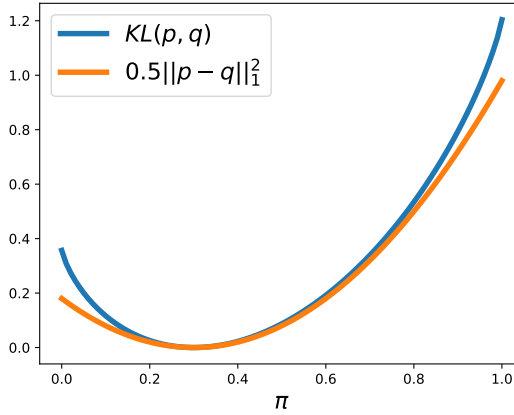


Figure 15.1: Graphical verification of Pinsker's inequality, $\frac{1}{2}\|\mathbf{p} - \mathbf{q}\|_1^2 \leq \text{KL}(\mathbf{p}, \mathbf{q})$, with $\mathbf{p} := (\pi, 1 - \pi)$ and $\mathbf{q} := (0.3, 0.7)$.

To see this, consider the function $F : t \rightarrow L(\mathbf{w}^* - t\nabla L(\mathbf{w}^*))$. If $\nabla L(\mathbf{w}^*) \neq 0$, then $F'(0) = -\|\nabla L(\mathbf{w}^*)\|_2^2 \neq 0$. Therefore, there exists a small $t > 0$ such that $F(t) < F(0)$, i.e., $L(\mathbf{w}^*)$ is not the minimum. However, if the function is not convex, the converse is a priori not true: finding a point that has a null gradient does not ensure that we have found a global minimum as illustrated in Figure 15.3.

For non-convex functions, a point with null gradient is called a **stationary point**. A stationary point may define a **local maximum** or a **local minimum**. Formally, $\hat{\mathbf{w}}$ is a local minimum if

$$\exists r > 0, \text{ s.t. } \forall \mathbf{v} \in \mathcal{W} \text{ satisfying } \|\mathbf{v} - \hat{\mathbf{w}}\| \leq r, \text{ we have } L(\mathbf{v}) \geq L(\hat{\mathbf{w}}).$$

A local maximum is defined similarly, except that $L(\mathbf{v}) \leq L(\hat{\mathbf{w}})$ in a neighborhood of $\hat{\mathbf{w}}$. For non-convex functions, convergence rates are therefore generally expressed in terms of convergence of the norm of the gradient $\|\nabla f(\mathbf{w}^t)\|_2$ towards 0. Such theoretical results do not ensure convergence to the global minimum but rather convergence to a point where no further progress may a priori be possible with just gradient information.

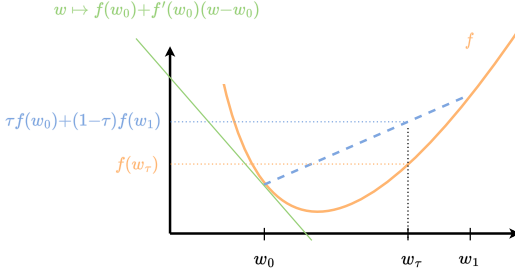


Figure 15.2: Convex function: any secant is above the function, any tangent is below the function, a point with zero gradient is a minimum.

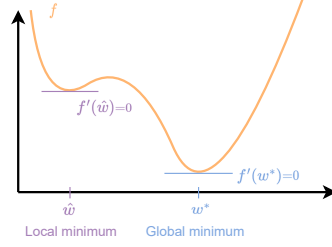


Figure 15.3: Non-convex function: a point with zero gradient is not necessarily the global minimum.

15.5 Performance guarantees

For a given class of functions, we can define the performance of an algorithm as the number of iterations the algorithm would need to find an ε -accurate solution as in Eq. (15.1). This is called the **computational complexity** of the algorithm, denoted

$$t = T(\varepsilon).$$

Alternatively, the performance of an algorithm can be stated in terms of **convergence rate**, i.e., the accuracy that the algorithm reaches after t iterations,

$$\varepsilon = R(t),$$

where R is a decreasing positive function vanishing as $t \rightarrow +\infty$. Usually, R incorporates properties of the function minimized, such as its smoothness constant β and information on the initial point, such as its function value. The corresponding computational complexity $T(\varepsilon)$ is then given as the minimum number of iterations t such that $R(t) \leq \varepsilon$,

$$T(\varepsilon) = \min\{t \in \mathbb{N} : R(t) \leq \varepsilon\}.$$

Convergence rates can generally be classified by considering the progress ratio on iteration t , defined by

$$\rho_t := \frac{R(t)}{R(t-1)}.$$

The asymptotic convergence rate is then defined by

$$\rho_\infty := \lim_{t \rightarrow +\infty} \rho_t.$$

We can classify the rates as follows.

1. **Sublinear convergence rates**, $\rho_\infty = 1$: the longer the algorithm runs, the slower it makes progress. That is, the relative progress eventually tends to stall as $t \rightarrow +\infty$. Examples of $R(t)$ in this category include $O(1/t)$, $O(1/t^2)$ or more generally $O(1/t^\alpha)$ for some $\alpha > 0$. This is equivalent to $T(\varepsilon) = O(\varepsilon^{-1/\alpha})$.
2. **Linear convergence rates**, $\rho_\infty = c \in (0, 1)$: the algorithm eventually reaches a state of constant relative progress at each iteration, leading to an overall rate $R(t) = O(\exp(-ct))$ for c depending on the properties of the objective. This corresponds to $T(\varepsilon) = O(c^{-1} \ln \varepsilon^{-1})$.
3. **Superlinear convergence rates**, $\rho_\infty = 0$: the relative progress is better at each new iteration. This can happen for, e.g., $R(t) = O(\exp(-t^2))$, leading to $T(\varepsilon) = O(\sqrt{\ln \varepsilon^{-1}})$ or $R(t) = O(\exp(-\exp(t)))$, also called a quadratic rate, leading to $T(\varepsilon) = O(\ln \ln \varepsilon^{-1})$.

This is illustrated in Fig. 15.4.

Note that the term “linear” may be misleading as the rates are in fact exponential. They are called “linear” because of their behavior in log scale.

Upper and lower bounds

The best performance of a class of algorithms equipped with a given oracle (e.g. first-order oracle) can be upper-bounded or lower-bounded. This allows to show that an algorithm with access limited to a certain type of oracle cannot theoretically do better than a certain number. For example, the computational complexity to minimize β -smooth functions restricted on $[0, 1]^P$ with first-order oracles is lower bounded by $\frac{c}{\varepsilon^P}$ (Nemirovski and Yudin, 1983, p. 1.1.7). For example, with $P = 10$ and $\varepsilon = 10^{-3}$, this gives 10^{30} iterations. Note that these results are pessimistic by construction. The actual performance of an algorithm on a specific instance of this function class may be much better than

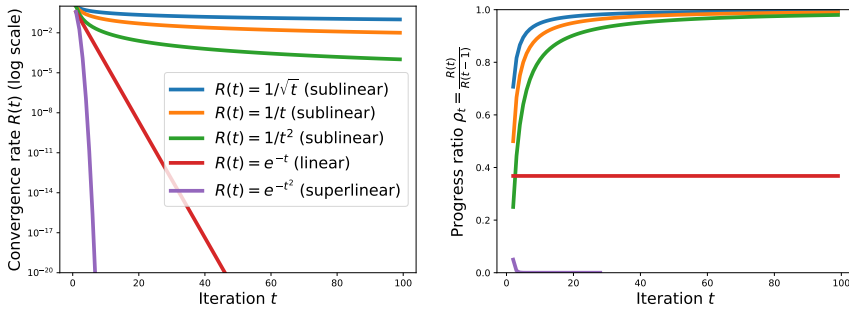


Figure 15.4: Left: convergence rates. **Right:** progress ratios. An algorithm with sublinear convergence rates eventually stops making progress. An algorithm with linear convergence rate eventually reaches a state of constant progress. An algorithm with superlinear convergence rate makes faster progress after each iteration.

this worst-case scenario, as it is the case with popular algorithms such as quasi-Newton methods. Better computational complexities can be achieved by further restricting the class of functions to the set of convex functions, which play a central role in optimization and many other fields.

Zero-order vs. first-order

For the class of smooth strongly convex functions, the computational complexity of the best first-order algorithm is (up to constant and logarithmic factors) P times better than that of the best zero-order algorithm (Nesterov, 2018; Nesterov and Spokoiny, 2017). This theoretical comparison shows that, while zero-order optimization algorithms may perform on par with first-order optimization algorithms for problems with a low dimension P , they can be much slower for high dimensional problems, i.e., $P \gg 1$.

In different settings, for example with stochastic oracles (Duchi *et al.*, 2015) or for different classes of functions, slightly different comparisons may be achieved, such as a \sqrt{P} factor instead of P . However, the same conclusion holds in the current frameworks considered: first-order optimization algorithms can provide fast rates that are dimension independent while the rates of zero-order optimization algorithms gen-

erally depend on the dimension of the problem, making them unfit for high-dimensional problems.

This explains the immense success of first-order algorithms for training neural networks. Fortunately, using reverse-mode autodiff, as studied in Chapter 8, it can be shown that computing a gradient has roughly the same complexity as evaluating the function itself Section 8.3.3

15.6 Summary

- The information available to us on a function can be formalized by the notion of **oracle**. Zero-order oracles can only evaluate the function; first-order oracles can also compute the gradient; second-order oracles can also compute the Hessian or the Hessian-vector product (HVP).
- Most optimization algorithms reviewed in this book can be viewed from a **variational perspective**, in which the next iteration is produced by optimizing a trade-off between an approximation of the function and a proximity term. Different approximations and different proximity terms lead to different algorithms.
- We also reviewed different classes of functions, and performance guarantees.

16

First-order optimization

16.1 Gradient descent

Gradient descent is one of the simplest algorithms in our toolbox to minimize a function. At each iteration, it moves along the negative gradient direction, scaled by a stepsize γ :

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \gamma \nabla L(\mathbf{w}^t). \quad (16.1)$$

The path taken by a gradient descent on a simple quadratic is illustrated in Fig. 16.1 for different choices of the stepsize.

16.1.1 Variational perspective

Consider the linear approximation of $L(\mathbf{w})$ around \mathbf{w}^t ,

$$L(\mathbf{w}) \approx L(\mathbf{w}^t) + \langle \nabla L(\mathbf{w}^t), \mathbf{w} - \mathbf{w}^t \rangle.$$

One can easily check that the gradient descent update in Eq. (16.1) can be rewritten as the solution of a minimization problem, namely,

$$\mathbf{w}^{t+1} = \arg \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}^t) + \langle \nabla L(\mathbf{w}^t), \mathbf{w} - \mathbf{w}^t \rangle + \frac{1}{2\gamma} \|\mathbf{w} - \mathbf{w}^t\|_2^2. \quad (16.2)$$

In words, a gradient descent update optimizes a trade-off between staying close to the current \mathbf{w}^t , thanks to the proximity term $\frac{1}{2\gamma} \|\mathbf{w} - \mathbf{w}^t\|_2^2$, and

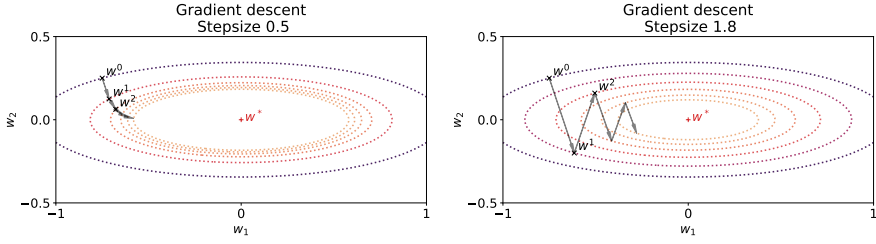


Figure 16.1: Trajectory taken by gradient descent on the objective $f(w) := 0.05w_1^2 + 0.5w_2^2$ with a small (left) or large (right) stepsize. In each case, the iterates follow the normal vectors to the contour lines (dashed lines): the negative gradients. A small stepsize leads to slow convergence but a larger stepsize induces oscillations.

minimizing the linearization of L around \mathbf{w}^t . Intuitively, by choosing γ sufficiently small, we ensure that the minimizer of the regularized linear approximation stays in a neighborhood where the linear approximation is valid. This viewpoint is useful to motivate gradient descent extensions.

16.1.2 Convergence for smooth functions

As long as $\nabla L(\mathbf{w}^t) \neq \mathbf{0}$, the function $L_t(\gamma) := L(\mathbf{w}^t - \gamma \nabla L(\mathbf{w}^t))$ has a negative derivative at 0, i.e., $L'_t(0) = -\|\nabla L(\mathbf{w}^t)\|_2^2$. Hence, as long as $\nabla L(\mathbf{w}^t) \neq \mathbf{0}$, there exists a stepsize ensuring a decrease in objective values at each iterate. However, without further assumptions, such a stepsize may depend on each iterate and may be infinitesimally small. To quantify the convergence of gradient descent with a constant stepsize, we restrict to the class of smooth functions. By applying Proposition 15.1 on the iterate of gradient descent, we obtain that

$$L(\mathbf{w}^{t+1}) \leq L(\mathbf{w}^t) - \gamma \|\nabla L(\mathbf{w}^t)\|_2^2 + \frac{\beta \gamma^2}{2} \|\nabla L(\mathbf{w}^t)\|_2^2.$$

Therefore, for β -smooth functions, by selecting $\gamma \leq \frac{1}{\beta}$, we get that

$$L(\mathbf{w}^{t+1}) - L(\mathbf{w}^t) \leq -\frac{\gamma}{2} \|\nabla L(\mathbf{w}^t)\|_2^2,$$

which illustrates the main mechanism behind gradient descent: each iteration decreases the objective by a constant times the norm of the gradient of the current iterate. This equation can further be summed

over all iterates up to T . This telescopes the objective values, leading to

$$\begin{aligned} \min_{t \in \{0, \dots, T-1\}} \|\nabla L(\mathbf{w}^t)\|_2^2 &\leq \frac{1}{T} \sum_{t=0}^{T-1} \|\nabla L(\mathbf{w}^t)\|_2^2 \\ &\leq \frac{2}{\gamma T} (L(\mathbf{w}^0) - L(\mathbf{w}^T)) \\ &\leq \frac{2}{\gamma T} (L(\mathbf{w}^0) - L^*), \end{aligned}$$

where we recall that L^* is the infimum of L . Therefore, after sufficiently many iterations, gradient descent finds a point whose gradient norm is arbitrarily small.

Non-convex case

Without further assumptions, i.e., in the non-convex case, the above result (i.e., convergence to a stationary point, measured by the gradient norm) is the best we may get in theory. Denoting $T_s(\varepsilon)$ the number of iterations needed for a gradient descent to output a point that is ε -stationary, i.e., $\|\nabla L(\hat{\mathbf{w}})\|_2 \leq \varepsilon$, we have $T_s(\varepsilon) \leq O(\varepsilon^{-2})$.

Convex case

By adding a convexity assumption on the objective, we can use the lower bound provided by the convexity assumption to ensure convergence to a minimum. Namely, for a β -smooth and convex function f , and with stepsize $\gamma \leq 1/\beta$, we have that (Nesterov, 2018)

$$L(\mathbf{w}^T) - L^* \leq \frac{1}{\gamma T} \|\mathbf{w}^0 - \mathbf{w}^*\|_2^2.$$

That is, we get a sublinear convergence rate, and the associated computational complexity to find a minimum is $T(\varepsilon) = O(1/\varepsilon)$.

Strongly convex case

If we further strengthen the assumptions by considering β -smooth, μ -strongly convex functions, the convergence rate of a gradient descent

can be shown to be (Nesterov, 2018), for any stepsize $\gamma \leq 1/\beta$,

$$\begin{aligned} L(\mathbf{w}^T) - L^\star &\leq (1 - \gamma\mu)^T \left(L(\mathbf{w}^0) - L^\star \right) \\ &\leq \exp(-\gamma\mu T) \left(L(\mathbf{w}^0) - L^\star \right). \end{aligned}$$

That is, we obtain a linear convergence rate and the associated computational complexity is $T(\varepsilon) = O(\ln \varepsilon^{-1})$. The above convergence rates may be further refined (Nesterov, 2018); we focused above on the simplest result for clarity.

Strong convexity can also be replaced by a weaker assumption, gradient-dominating property (Polyak, 1963), i.e., $\|\nabla L(\mathbf{v})\|_2^2 \geq c(L(\mathbf{v}) - L^\star)$ for some constant c and any $\mathbf{v} \in \mathcal{W}$. A convex, gradient-dominating function can also be minimized at a linear rate.

16.1.3 Momentum and accelerated variants

We started with gradient descent as a simple example of first-order optimization algorithm. However, different optimization algorithms can be designed from the access to first-order oracles and the knowledge of the class of functions considered. For example, consider quadratic convex functions $\mathbf{w} \mapsto \frac{1}{2}\mathbf{w}^\top \mathbf{A}\mathbf{w} + \mathbf{b}^\top \mathbf{w}$, that are a basic example of smooth strongly convex functions if \mathbf{A} is positive definite. An optimal method in this case is the heavy-ball method of Polyak (1964), that can be written as

$$\begin{aligned} \mathbf{v}^{t+1} &:= \nu \mathbf{v}^t - \gamma \nabla L(\mathbf{w}^t) \\ \mathbf{w}^{t+1} &:= \mathbf{w}^t + \mathbf{v}^{t+1}. \end{aligned}$$

The heavy-ball method uses an additional variable \mathbf{v}^t , that can be interpreted as the velocity of a ball driven by the negative gradient to converge towards a minimum. Intuitively, this additional velocity circumvents the oscillations that a gradient descent may present as illustrated in Fig. 16.2 compared to Fig. 16.1. For $\nu = 0$, we recover usual gradient descent. For $\nu > 0$, the velocities accumulate a form of an inertia momentum, where ν is interpreted as the “mass” of the ball. In terms of convergence rates, the heavy-ball method can be shown to converge linearly similarly to gradient descent, but with a rate

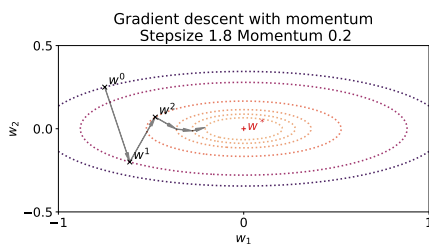


Figure 16.2: Trajectory taken by gradient descent with momentum. Compared to gradient descent without momentum, for the same stepsize, the oscillations previously observed in Fig. 16.1 are no longer present, and the algorithm converges faster to the minimum.

$O(\exp(-T\sqrt{\mu/\beta}))$ for appropriate choices of ν, γ . In comparison, by choosing an optimal stepsize for the gradient descent, its convergence rate is $O(\exp(-T\mu/\beta))$ which is provably worse, as we always have $\mu/\beta \leq 1$.

Beyond the case of quadratic functions, accelerated variants of gradient descent for convex or strongly convex functions have been developed by Nesterov (2018). Such variants have inspired the design of optimization algorithms in stochastic settings presented below.

16.2 Stochastic gradient descent

In machine learning, we are usually interested in minimizing the **expected loss** of the model over the data distribution ρ :

$$\min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}) := \mathbb{E}_{S \sim \rho} [L(\mathbf{w}; S)].$$

For example, L is often set to $L(\mathbf{w}; S) := \ell(Y, f(X, \mathbf{w}))$, where ℓ is a loss function, f is a neural network and $S = (X, Y)$ is a random pair, composed of an input X and an associated target Y , sampled from ρ . In this setting, since the data distribution ρ is generally unknown and may be infinite, we cannot exactly evaluate the expected loss $L(\mathbf{w})$ or its gradient $\nabla L(\mathbf{w})$.

In practice, we are often given a fixed dataset of n pairs $\mathbf{s}_i = (\mathbf{x}_i, \mathbf{y}_i)$. This is a special case of the expected loss setting, since this can be seen

as a empirical distribution $\rho = \rho_n$

$$L(\mathbf{w}) = \mathbb{E}_{S \sim \rho_n} [L(\mathbf{w}; S)] = \frac{1}{n} \sum_{i=1}^n L(\mathbf{w}; (X_i, Y_i)).$$

The gradient of $L(\mathbf{w})$ is then

$$\nabla L(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \nabla L(\mathbf{w}; (\mathbf{x}_i, \mathbf{y}_i)).$$

In this case, we see that the **full** gradient $\nabla L(\mathbf{w})$, as needed by gradient descent, is the average of the **individual** gradients. That is, the cost of computing $\nabla L(\mathbf{w})$ is proportional to the number of training points n . For n very large, that is a very large amount of samples, this computational cost can be prohibitive. Stochastic gradients circumvent this issue.

16.2.1 Stochastic gradients

Usually, even if we do not know ρ , we can sample from it, i.e., we have access to samples $S \sim \rho$. We can then use a **stochastic gradient** of the form $\nabla L(\mathbf{w}; S)$ as a random estimate of $\nabla L(\mathbf{w})$. This may look like a rough estimate but, on average, this is a valid approximation since

$$\mathbb{E}_{S \sim \rho} [\nabla L(\mathbf{w}; S)] = \nabla L(\mathbf{w}).$$

We say that $\nabla L(\mathbf{w}; S)$ is an **unbiased estimator** of $\nabla L(\mathbf{w})$. To further improve the approximation, we may also consider **mini-batch** estimates by sampling $m \ll n$ data points $S_i := (X_i, Y_i)$ and using $\frac{1}{m} \sum_{i=1}^m \nabla L(\mathbf{w}; S_i)$, whose expectation still matches $\nabla L(\mathbf{w})$, while potentially reducing the approximation error by averaging multiple stochastic gradients. Computationally, the main advantage is that the cost is now proportional to m instead of n .

In whole generality, one can consider stochastic first-order oracles defined below.

Definition 16.1 (Stochastic first-order oracles). A **stochastic first-order oracle** of an expected objective $L(\mathbf{w})$ is a random estimate $g(\mathbf{w}; S)$ of $\nabla L(\mathbf{w})$ with S sampled according to some distribution

q . A stochastic gradient is said to be an **unbiased estimator** if

$$\mathbb{E}_{S \sim q} [g(\mathbf{w}; S)] = \nabla L(\mathbf{w}).$$

The **variance** of a stochastic gradient is

$$\mathbb{E}_{S \sim q} [\|g(\mathbf{w}; S) - \nabla L(\mathbf{w})\|_2^2].$$

When $q = \rho$, we recover stochastic gradients. When q is the product of m independent samples according to p , we recover mini-batch stochastic gradients. First-order stochastic optimization algorithms build upon stochastic first-order oracles to approximately find the minimum of the expected objective. In such a setting, the iterates of the algorithm are by definition random. Convergence rates therefore need to be expressed in probabilistic terms by considering for example the expected objective value according to the randomness of the oracles.

16.2.2 Vanilla SGD

Equipped with a stochastic first-order oracle, such as (mini-batch) stochastic gradients, we can define **stochastic gradient descent** as

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \gamma g(\mathbf{w}^t; S^t) \quad \text{where } S^t \sim q.$$

We assume that S^t is independent of \mathbf{w}^t . Compared to the usual gradient descent, the main impediment of the stochastic setting is the additional noise induced by the stochastic estimates: their variance.

For example, consider applying a stochastic gradient descent on the expectation of β -smooth convex functions $L(\mathbf{w}; \mathbf{s})$ with unbiased oracles. To harness the randomness of the iterates, consider after T iterations outputting the average of the first T iterates, that is $\bar{\mathbf{w}}^T := \frac{1}{T} \sum_{t=1}^T \mathbf{w}^t$. Moreover, suppose that the variance of the stochastic first-order oracles is bounded by σ^2 for all minimizers \mathbf{w}^* of L . Denoting by $\mathbb{E}_{S_0, \dots, S_{T-1}}$ the randomness associated to the stochastic oracles, we have then that for a stepsize $\gamma \leq 1/(4\beta)$, (Lan, 2012),

$$\mathbb{E}_{S_0, \dots, S_{T-1}} [L(\bar{\mathbf{w}}^T)] - L^* \leq \frac{1}{\gamma T} \|\mathbf{w}_0 - \mathbf{w}^*\|_2^2 + 2\gamma\sigma^2.$$

The resulting convergence rate illustrates that a stochastic gradient descent converges to the minimum of the expected objective up to a

constant term depending on the variance of the oracle and the stepsize. One can diminish the variance by considering mini-batches: if the variance of a single stochastic gradient is σ_1^2 , considering a mini-batch of m gradients reduces the variance of the corresponding oracle to $\sigma_m = \sigma_1^2/m$. To decrease the additional term, one may also decrease the stepsizes over the iterations. For example, by choosing a decreasing stepsize like $\gamma^t = t^{-1/2}$, the convergence rate is then of the order $O((\|\mathbf{w}_0 - \mathbf{w}^*\|_2^2 + \sigma^2 \ln t)/\sqrt{t})$. The stepsize can also be selected as a constant γ_0 that decreases the average objective for the first T_0 iterations and reduced by a multiplicative factor at regular intervals like $\gamma_j = \rho\gamma_{j-1}$ for $\rho \in (0, 1)$ to handle iterations between T_j, T_{j+1} . Alternative stepsize schedules such as a cosine decay (Loshchilov and Hutter, 2016) have recently become popular.

The literature on alternative optimization schemes for stochastic optimization is still rapidly evolving, with new heuristics regularly proposed. We present below two popular techniques.

16.2.3 Momentum variants

Accelerated optimization algorithms developed in the deterministic setting may be extended to the stochastic setting. For example, the heavy-ball method can be adapted to the stochastic setting, leading to stochastic gradient descent with **momentum** (Sutskever *et al.*, 2013) generally implemented as

$$\begin{aligned}\mathbf{v}^{t+1} &:= \nu \mathbf{v}^t + \mathbf{g}(\mathbf{w}^t; S^t) \\ \mathbf{w}^{t+1} &:= \mathbf{w}^t - \gamma \mathbf{v}^{t+1}.\end{aligned}$$

As mentioned earlier the momentum method can be modified to handle non-quadratic smooth strongly convex functions. This leads to Nesterov's accelerated method in the deterministic setting. This has been adapted to the stochastic with a so-called **Nesterov momentum** (Sutskever *et al.*, 2013)

$$\begin{aligned}\mathbf{v}^{t+1} &:= \nu \mathbf{v}^t + \mathbf{g}(\mathbf{w}^t + \nu \mathbf{v}^t; S^t) \\ \mathbf{w}^{t+1} &:= \mathbf{w}^t - \gamma \mathbf{v}^{t+1}.\end{aligned}$$

16.2.4 Adaptive variants

In any gradient descent-like algorithm, selecting the stepsize is key for good performance. While a constant stepsize may be used if the function is smooth, we may not know in advance the smoothness constant of the objective, which means that additional procedures may be required to select appropriately the stepsize. In the deterministic case, line-searches such as the Armijo or Wolfe's rules (Wright and Nocedal, 1999) can be used to check whether the selected stepsize decreases sufficiently the objective at each iteration. Such rules have been adapted in the stochastic setting (Vaswani *et al.*, 2019).

Another way to decrease the sensitivity of the algorithm with respect to the stepsize has been to estimate first and second-order moments of the gradients and use the latter as a form of preconditioning to smooth the trajectory of the iterates. This led to the popular **Adam** optimizer (Kingma and Ba, 2014). It takes the form,

$$\begin{aligned}\mathbf{m}^{t+1} &:= \nu_1 \mathbf{m}^t + (1 - \nu_1) \mathbf{g}^t \\ \mathbf{v}^{t+1} &:= \nu_2 \mathbf{v}^t + (1 - \nu_2) (\mathbf{g}^t)^2 \\ \hat{\mathbf{m}}^{t+1} &:= \mathbf{m}^{t+1} / (1 - \nu_1^t) \\ \hat{\mathbf{v}}^{t+1} &:= \mathbf{v}^{t+1} / (1 - \nu_2^t) \\ \mathbf{w}^{t+1} &:= \mathbf{w}^t - \gamma \hat{\mathbf{m}}^{t+1} / \left(\sqrt{\hat{\mathbf{v}}^{t+1} + \varepsilon} \right),\end{aligned}$$

where $\mathbf{g}^t := \mathbf{g}(\mathbf{w}^t; S^t)$, $(\mathbf{g}^t)^2$ denotes the element-wise square of \mathbf{g}^t and $\nu_1, \nu_2, \gamma, \varepsilon$ are hyper-parameters of the algorithm. Numerous variants exist, such as varying the stepsize γ above along the iterations.

16.3 Projected gradient descent

Oftentimes, we seek to find the solution of a minimization problem subject to **constraints** on the variables, of the form

$$\min_{\mathbf{w} \in \mathcal{C}} L(\mathbf{w}), \tag{16.3}$$

where $\mathcal{C} \subseteq \mathcal{W} = \mathbb{R}^P$ is a set of constraints. We say that an approximate solution $\hat{\mathbf{w}}$ to Eq. (16.3) is **feasible** if $\hat{\mathbf{w}} \in \mathcal{C}$. Naturally, the design

of algorithms for the constrained setting now depends, not only on information about L , but also on information about \mathcal{C} .

Similarly to L , different **oracles** can be considered about \mathcal{C} . One of the most commonly used oracle is the **Euclidean projection**

Definition 16.2 (Euclidean projection). The Euclidean projection onto the set \mathcal{C} is defined by

$$\text{proj}_{\mathcal{C}}(\mathbf{w}) := \arg \min_{\mathbf{v} \in \mathcal{C}} \|\mathbf{w} - \mathbf{v}\|_2^2.$$

This projection, which is well-defined when \mathcal{C} is a convex set, can be used in projected gradient descent, that we briefly review below. Typically, the projection on a particular set \mathcal{C} requires a dedicated algorithm to compute it.

Other possible oracles are **linear maximization oracles** (LMO) used in Frank-Wolfe algorithms and **Bregman projection oracles**, used in mirror descent algorithms. The algorithm choice can be dictated by what oracle about \mathcal{C} is available.

16.3.1 Variational perspective

Projected gradient descent is a natural generalization of gradient descent, based on the Euclidean projection oracle. Its iterates read

$$\mathbf{w}^{t+1} := \text{proj}_{\mathcal{C}}(\mathbf{w}^t - \gamma \nabla L(\mathbf{w}^t)).$$

At each iteration, we attempt to decrease the objective by moving along the negative gradient direction, while ensuring that the next iterate remains feasible, thanks to the projection step.

Similarly to the variational perspective of gradient descent in Eq. (16.2), the projected gradient descent update is equivalent to

$$\mathbf{w}^{t+1} = \arg \min_{\mathbf{w} \in \mathcal{C}} L(\mathbf{w}^t) + \langle \nabla L(\mathbf{w}^t), \mathbf{w} - \mathbf{w}^t \rangle + \frac{1}{2\gamma} \|\mathbf{w} - \mathbf{w}^t\|_2^2.$$

This shows that projected gradient descent minimizes a trade-off between staying close to \mathbf{w}^t and minimizing the linearization of L around \mathbf{w}^t , while staying in \mathcal{C} .

In terms of convergence rates, they remain the same as gradient descent (Nesterov, 2018). For example, projected gradient descent on a smooth convex function still converges at a rate $R(T) = O(1/T)$.

There are numerous extensions of vanilla projected gradient descent. Similarly to gradient descent, the stepsize can be automatically adjusted using linesearch techniques and there exists accelerated variants. If we replace $\nabla L(\mathbf{w})$ with a stochastic gradient $\nabla L(\mathbf{w}; S)$, we obtain a stochastic projected gradient descent.

16.3.2 Optimality conditions

In the unconstrained case, a minimum necessarily has a zero gradient. In the constrained setting, there may not be any feasible parameters with zero gradient. Instead, the optimality of a point is characterized by the fact that no better solution can be found by moving along the gradient at that point, while staying in the constraints. Formally, it means that for any $\gamma > 0$, a minimizer \mathbf{w}^* of L on \mathcal{C} satisfies

$$\mathbf{w}^* = \text{proj}_{\mathcal{C}}(\mathbf{w}^* - \gamma \nabla L(\mathbf{w}^*)).$$

It can be shown that this condition is equivalent (Nesterov, 2018) to

$$\langle \nabla L(\mathbf{w}^*), \mathbf{w} - \mathbf{w}^* \rangle \geq 0 \quad \forall \mathbf{w} \in \mathcal{C}.$$

16.3.3 Commonly-used projections

We now briefly review a few useful Euclidean projections.

- If $\mathcal{C} = \mathbb{R}^P$, we obviously have

$$\text{proj}_{\mathcal{C}}(\mathbf{w}) = \mathbf{w}.$$

Therefore, in the unconstrained setting, projected gradient descent indeed recovers gradient descent.

- If $\mathcal{C} = [a, b]^P$ (box constraints), we have

$$\text{proj}_{\mathcal{C}}(\mathbf{w}) = \text{clip}(\mathbf{w}, a, b) := \min\{\max\{\mathbf{w}, a\}, b\}.$$

where the min and max are applied coordinate-wise.

- As a special case of the above, if $\mathcal{C} = \mathbb{R}_+^P$ (non-negative orthant),

$$\text{proj}_{\mathcal{C}}(\mathbf{w}) = \max\{\mathbf{w}, 0\},$$

also known as non-negative part or ReLu.

- If $\mathcal{C} = \Delta^P$ (unit probability simplex),

$$\text{proj}_{\mathcal{C}}(\mathbf{w}) = \max\{\mathbf{w} - \tau \mathbf{1}, 0\},$$

where $\tau \in \mathbb{R}$ is a constant ensuring that $\text{proj}_{\mathcal{C}}(\mathbf{w})$ normalizes to 1. It is known that τ can be found in $O(P \log P)$ using a sort. This can be improved to $O(P)$ using a median-finding like algorithm.

16.4 Proximal gradient method

The constrained setting (with \mathcal{C} a convex set) can be recast as unconstrained optimization, by extending our analysis to functions taking infinite values. Let us denote the indicator function of the set \mathcal{C} by

$$\iota_{\mathcal{C}}(\mathbf{w}) := \begin{cases} 0 & \text{if } \mathbf{w} \in \mathcal{C} \\ +\infty & \text{otherwise} \end{cases}.$$

Clearly, the constrained problem in Eq. (16.3) can then be rewritten as

$$\min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}) + \iota_{\mathcal{C}}(\mathbf{w}).$$

This suggests that constrained optimization is a special case of **composite objectives** of the form

$$\min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}) + \Omega(\mathbf{w}),$$

where Ω is a convex but potentially non-differentiable function. We assume that we have access to an oracle associated with Ω called the **proximal operator**.

Definition 16.3 (Proximal operator). The proximal operator associated with $\Omega: \mathcal{W} \rightarrow \mathbb{R}$ is

$$\text{prox}_{\Omega}(\mathbf{w}) := \arg \min_{\mathbf{v} \in \mathcal{W}} \frac{1}{2} \|\mathbf{w} - \mathbf{v}\|_2^2 + \Omega(\mathbf{v}).$$

This leads to the proximal gradient method, reviewed below.

16.4.1 Variational perspective

With this method, the update reads

$$\mathbf{w}^{t+1} = \text{prox}_{\gamma\Omega}(\mathbf{w}^t - \gamma\nabla L(\mathbf{w}^t)).$$

This update again enjoys an intuitive variational perspective, namely,

$$\mathbf{w}^{t+1} = \arg \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}^t) + \langle \nabla L(\mathbf{w}^t), \mathbf{w} - \mathbf{w}^t \rangle + \frac{1}{2\gamma} \|\mathbf{w} - \mathbf{w}^t\|_2^2 + \Omega(\mathbf{w}).$$

That is, we linearize L around \mathbf{w}^t , but keep Ω as is.

The proximal gradient method is popularly used when the objective function contains a sparsity-inducing regularizer Ω . For example, for the LASSO (Tibshirani, 1996), which aims at predicting targets $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^N$ from observations $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top \in \mathbb{R}^{N \times P}$, we set $L(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$ and $\Omega(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$, where $\lambda > 0$ controls the regularization strength. In this case, prox_Ω is the so-called soft-thresholding operator (see below).

Convergence guarantees of the proximal gradient method remain the same as for gradient descent, such as a $O(1/T)$ rate for smooth convex functions.

16.4.2 Optimality conditions

An optimal solution of the problem is characterized by the **fixed point** equation

$$\mathbf{w}^* = \text{prox}_{\gamma\Omega}(\mathbf{w}^* - \gamma\nabla L(\mathbf{w}^*)),$$

for all $\gamma > 0$ (Nesterov, 2018). In other words, the proximal gradient method (which includes gradient descent and projected gradient descent as special cases), can be seen as fixed point iteration schemes. Such a viewpoint suggests using acceleration methods from the fixed point literature such as Anderson acceleration (Pollock and Rebholz, 2021). It is also useful when designing implicit differentiation schemes as presented in Chapter 8.

16.4.3 Commonly-used proximal operators

We now briefly review a few useful proximal operators.

- If $\Omega(\mathbf{w}) = 0$, we have

$$\text{prox}_{\gamma\Omega}(\mathbf{w}) = \mathbf{w}.$$

Therefore, with this proximal operator, the proximal gradient method recovers gradient descent.

- If $\Omega(\mathbf{w}) = \iota_C(\mathbf{w})$, we have

$$\text{prox}_{\gamma\Omega}(\mathbf{w}) = \text{proj}_C(\mathbf{w}).$$

Therefore, with this proximal operator, the proximal gradient method recovers projected gradient descent.

- If $\Omega(\mathbf{w}) = \lambda\|\mathbf{w}\|_1$, we have

$$\text{prox}_{\gamma\Omega}(\mathbf{w}) = (\text{sign}(\mathbf{w}) \cdot \max(|\mathbf{w}| - \gamma\lambda, 0)),$$

where the operations are applied coordinate-wise. This is the so-called soft-thresholding operator.

- $\Omega(\mathbf{w}) = \lambda \sum_{g \in G} \|\mathbf{w}_g\|_2$ where G is a partition of $[P]$ and \mathbf{w}_g denotes the subvector restricted to g , then we have

$$\left[\text{prox}_{\gamma\Omega}(\mathbf{w}) \right]_g = \max(1 - \lambda \cdot \gamma / \|\mathbf{w}_g\|_2, 0) \mathbf{w}_g,$$

which is used in the group lasso (Yuan and Lin, 2006) and can be used to encourage group sparsity.

For a review of more proximal operators, see for instance (Bach *et al.*, 2012; Parikh, Boyd, *et al.*, 2014).

16.5 Summary

- From a variational perspective, gradient descent is the algorithm obtained when linearizing the objective function and using a quadratic regularization term.

- Projected gradient descent is the algorithm obtained when there is an additional constraint (the Euclidean projection naturally appearing, due to the quadratic regularization term).
- When the objective is the sum of a differentiable function and a non-differentiable function, proximal gradient is the algorithm obtained when the differentiable function is linearized but the non-differentiable function is kept as is.
- We also reviewed various stochastic gradient based algorithms, including vanilla SGD, SGD with momentum and Adam.

17

Second-order optimization

We review in this chapter methods whose iterations take the form

$$\mathbf{w}^{t+1} := \mathbf{w}^t - \gamma^t B^t \nabla L(\mathbf{w}^t),$$

where γ^t is a stepsize and B^t is a pre-conditioning matrix involving second-order derivatives.

17.1 Newton's method

17.1.1 Variational perspective

We saw in Eq. (16.2) that gradient descent can be motivated from a variational perspective, in which we use a linear approximation of the objective around the current iterate, obtained from the current gradient. Similarly, if we have access not only to the gradient but also to the Hessian of the objective, we can use a quadratic approximation of the objective around the current iterate. More precisely, given a function $L(\mathbf{w})$, we may consider minimizing the second-order Taylor approximation of $L(\mathbf{w})$ around the current iterate \mathbf{w}^t ,

$$L(\mathbf{w}) \approx L(\mathbf{w}^t) + \langle \nabla L(\mathbf{w}^t), \mathbf{w} - \mathbf{w}^t \rangle + \frac{1}{2} \langle \mathbf{w} - \mathbf{w}^t, \nabla^2 L(\mathbf{w}^t) (\mathbf{w} - \mathbf{w}^t) \rangle.$$

Newton's method simply iteratively minimizes this quadratic approximation around the current iteration \mathbf{w}^t , namely,

$$\mathbf{w}^{t+1} = \arg \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}^t) + \langle \nabla L(\mathbf{w}^t), \mathbf{w} - \mathbf{w}^t \rangle + \frac{1}{2} \langle \mathbf{w} - \mathbf{w}^t, \nabla^2 L(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) \rangle. \quad (17.1)$$

If the Hessian is positive definite at \mathbf{w}^t , which we denote by $\nabla^2 L(\mathbf{w}^t) \succ 0$, then the minimum is well-defined and unique (this is for example the case if L is strictly convex). The iterates can then be written analytically as

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \nabla^2 L(\mathbf{w}^t)^{-1} \nabla L(\mathbf{w}^t).$$

If the Hessian is not positive definite, the minimum may not be defined. Ignoring this issue and taking the analytical formulation could be dangerous, as it could amount to computing the maximum of the quadratic instead if, for example, the quadratic was strictly concave (i.e., $\nabla^2 L(\mathbf{w}) \prec 0$).

17.1.2 Regularized Newton method

A simple technique to circumvent this issue consists in adding a regularization term to the Hessian. Namely, from a variational viewpoint, we can add a proximity term $\frac{1}{2} \|\mathbf{w} - \mathbf{w}^t\|_2^2$, encouraging to stay close to the current \mathbf{w}^t . The iterates of this regularized Newton method then take the form

$$\begin{aligned} \mathbf{w}^{t+1} = \arg \min_{\mathbf{w} \in \mathcal{W}} & L(\mathbf{w}^t) + \langle \nabla L(\mathbf{w}^t), \mathbf{w} - \mathbf{w}^t \rangle + \frac{1}{2} \langle \mathbf{w} - \mathbf{w}^t, \nabla^2 L(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) \rangle \\ & + \frac{\eta^t}{2} \|\mathbf{w} - \mathbf{w}^t\|_2^2, \end{aligned}$$

where η^t controls the regularization strength. Assuming $\eta^t > 0$ is strong enough to make $\nabla^2 L(\mathbf{w}^t) + \eta^t \mathbf{I}$ positive-definite, we have

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{d}^t,$$

where we defined the direction

$$\mathbf{d}^t := (\nabla^2 L(\mathbf{w}^t) + \eta^t \mathbf{I})^{-1} \nabla L(\mathbf{w}^t). \quad (17.2)$$

Other techniques to circumvent this issue include using cubic regularization and modifying the spectral decomposition of the Hessian, by thresholding the eigenvalues or taking their absolute values. We refer the interested reader to, e.g., (Nesterov, 2018; Wright and Nocedal, 1999) for more details.

17.1.3 Approximate direction

We observe a main impediment for implementing such a second-order optimization algorithm: even if we had access to the Hessian of the objective for free and this Hessian was positive definite, computing the exact direction \mathbf{d}^t in Eq. (17.2) requires computing an inverse-Hessian vector product (IHVP) with the gradient $\nabla L(\mathbf{w}^t)$. Doing so exactly requires solving a linear system

$$(\nabla^2 L(\mathbf{w}^t) + \eta^t \mathbf{I})\mathbf{d}^t = \nabla L(\mathbf{w}^t),$$

which a priori takes $O(P^3)$ time. In practice, however, we can compute IHVPs approximately, as explained in Section 9.4.

17.1.4 Convergence guarantees

While implementing Newton's method comes at a higher computational cost, it can also benefit from faster convergence rates. Briefly, if Newton's method is initialized at a point $\mathbf{w}^0 \in \mathcal{W}$ close enough from the minimizer \mathbf{w}^* of a μ -strongly convex function with M -Lipschitz continuous Hessian (namely $\|\mathbf{w}^0 - \mathbf{w}^*\|_2 \leq \frac{2\mu}{3M}$), then Newton's method converges at a quadratic rate (Nesterov, 2018), that is, $R(t) \leq O(\exp(\exp(-t)))$ (see Section 15.5 for a brief introduction to performance guarantees). This is far superior to gradient descent. Such an efficiency motivated the development of interior point methods, that have been a breakthrough in constrained optimization, thanks to the use of log-barrier penalties (Nesterov, 2018).

17.1.5 Linesearch

In practice, we may not have access to an initial point close enough from the minimizer. In that case, even for strictly convex functions for

which Newton's steps are well-defined, taking $\mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{d}^t$ may not ensure a decrease of the objective values. Nevertheless, the direction \mathbf{d}^t may define a descent direction as defined below.

Definition 17.1 (Descent direction). A point $\mathbf{d} \in \mathcal{W}$ defines a **descent direction** $-\mathbf{d}$ for an objective L at \mathbf{w} , if there exists a positive stepsize $\gamma > 0$ such that

$$L(\mathbf{w} - \gamma \mathbf{d}) \leq L(\mathbf{w}).$$

If L is differentiable, $-\mathbf{d}$ is a descent direction if $\langle -\mathbf{d}, \nabla L(\mathbf{w}) \rangle < 0$.

For Newton's method without regularization, $\mathbf{d}^t = \nabla^2 L(\mathbf{w}^t)^{-1} \nabla L(\mathbf{w}^t)$ is then a descent direction at \mathbf{w}^t , as long as $\nabla L(\mathbf{w}^t) \neq 0$ and $\nabla^2 L(\mathbf{w}^t) \succ 0$. If $\nabla^2 L(\mathbf{w}^t) \not\succ 0$, choosing $\eta^t > 0$ such that $\nabla^2 L(\mathbf{w}^t) + \eta^t \mathbf{I} \succ 0$, also ensures that $\mathbf{d}^t = -(\nabla^2 L(\mathbf{w}^t) + \eta^t \mathbf{I})^{-1} \nabla L(\mathbf{w}^t)$ is a descent direction (as long as $\nabla L(\mathbf{w}^t) \neq 0$). Newton's method is then generally equipped with a linesearch method that attempts to take steps of the form

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \gamma^t \mathbf{d}^t$$

with γ^t chosen as the largest stepsize among $\{\rho^\tau, \tau \in \mathbb{N}\}$ for $\rho \in (0, 1)$ until a sufficient decrease of the objective is satisfied such as, for $c \in (0, 1)$,

$$L(\mathbf{w}^t - \gamma^t \mathbf{d}^t) \leq L(\mathbf{w}^t) - c\gamma^t \langle \nabla L(\mathbf{w}^t), \nabla^2 L(\mathbf{w}^t)^{-1} \nabla L(\mathbf{w}^t) \rangle.$$

For strongly convex functions, such an implementation exhibits two phases: a first phase during which Newton's steps are “damped” by using a stepsize $\gamma^t < 1$ and a second phase of super-fast convergence during which stepsizes $\gamma^t = 1$ are taken, and the objective decreases very fast. Even far from the optimum, Newton directions can advantageously adapt to the local geometry of the objective to speed-up convergence compared to a regular gradient descent as explained below.

17.1.6 Geometric interpretation

To understand the efficiency of Newton's method compared to gradient descent, consider the minimization of a simple quadratic

$$L(\mathbf{w}) = \frac{1}{2}aw_1^2 + \frac{1}{2}bw_2^2$$

for $a \gg b \geq 0$, as illustrated in Fig. 17.1. A gradient descent moves along the directions $\nabla L(\mathbf{w}) = (aw_1, bw_2)^\top$ and its stepsize is limited by the variations in the first coordinate leading to some oscillations. If we were simply rescaling the gradient by (a, b) , i.e., taking steps of the form

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \gamma \mathbf{diag}(a^{-1}, b^{-1}) \nabla L(\mathbf{w}^t),$$

the variations in both coordinates would be normalized to one and the stepsize could simply be chosen to $\gamma = 1$ to directly get \mathbf{w}^* . In other words, by adapting the geometry of the directions with the geometry induced by the objective, we can circumvent the oscillations.

That's exactly what Newton's method does by modifying the gradient direction using the inverse of the Hessian. Formally, at iteration t , consider the modified objective

$$\tilde{L}(\mathbf{v}) = L(A\mathbf{v}) \text{ for } A = \nabla^2 L(\mathbf{w}^t)^{-1/2},$$

with L strictly convex and A the inverse matrix square root of the Hessian. One easily verifies that a Newton step is equivalent to a gradient step on \tilde{L} , that is,

$$\mathbf{v}^{t+1} = \mathbf{v}^t - \nabla \tilde{L}(\mathbf{v}^t) \iff \mathbf{w}^{t+1} = \mathbf{w}^t - (\nabla^2 L(\mathbf{w}^t))^{-1} \nabla L(\mathbf{w}^t)$$

where

$$\mathbf{w}^t = A\mathbf{v}^t = \nabla^2 L(\mathbf{w}^t)^{-1/2} \mathbf{v}^t.$$

In the geometry induced by A , the objective is generally better conditioned as illustrated in Fig. 17.1. This explains the efficiency of Newton's method. In particular for any strongly convex quadratic, a Newton step reaches the optimum in one iteration, while a gradient step can take many more iterations.

17.1.7 Stochastic Newton's method

Consider now an expected loss

$$\min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}) := \mathbb{E}_{S \sim \rho} [L(\mathbf{w}; S)].$$

In that case, an estimate of the Hessian can be constructed just like for the gradient using that

$$\mathbb{E}_{S \sim \rho} [\nabla^2 L(\mathbf{w}; S)] = \nabla^2 L(\mathbf{w}).$$

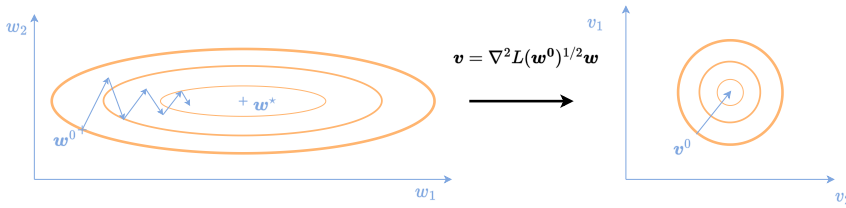


Figure 17.1: Left: Minimization of a quadratic $L(\mathbf{w}) := \frac{1}{2}a\mathbf{w}_1^2 + \frac{1}{2}b\mathbf{w}_2^2$ by gradient descent. For $a \gg b \geq 0$, gradient descent typically oscillates. Right: minimization by Newton's method amounts to change the geometry of the problem to avoid oscillations.

Denote then

$$g(\mathbf{w}; S) \approx \nabla L(\mathbf{w}), \quad H(\mathbf{w}; S') \approx \nabla^2 L(\mathbf{w})$$

some stochastic estimates of respectively of the gradient and the Hessian with S, S' independently drawn from p or from mini-batch approximations with varying mini-batch sizes. One implementation of a **stochastic Newton method** can then be

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \gamma^t (H(\mathbf{w}^t; S') + \eta^t \mathbf{I})^{-1} g(\mathbf{w}^t; S),$$

for $\eta^t \geq 0$ such that $(H(\mathbf{w}^t; S') + \eta^t \mathbf{I})^{-1} \succ 0$ and γ^t fixed or chosen to satisfy some sufficient decrease condition. We refer the interested reader to, e.g., (Xu *et al.*, 2020), for more details and variants.

17.2 Gauss-Newton method

Newton's method (17.1) is usually not properly defined for non-convex objective functions, since the Hessian may not be positive definite at the current iterate. We saw in Section 9.2 that the Gauss-Newton matrix can be used to define a positive-semidefinite approximation of the Hessian. Here, we revisit the Gauss-Newton method from a variational and **partial linearization** perspective. While the original Gauss-Newton method originates from nonlinear least-squares, we will first describe an extension to arbitrary convex loss functions, since it is both more general and easier to explain.

17.2.1 With exact outer function

Consider a composite objective of the form

$$L(\mathbf{w}) := \ell(f(\mathbf{w})),$$

where $\ell : \mathcal{M} \rightarrow \mathbb{R}$ is a **convex** function, such as a convex loss function applied on a given sample, and $f : \mathcal{W} \rightarrow \mathcal{M}$ is a **nonlinear** function, such as a neural network with parameters $\mathbf{w} \in \mathcal{W}$, evaluated on the same sample. We saw that gradient descent and Newton's method amount to using **linear** and **quadratic** approximations of $L(\mathbf{w})$ around the current iterate \mathbf{w}^t , respectively. As a middle ground between the two, the Gauss-Newton method uses the linearization of f around \mathbf{w}^t

$$f(\mathbf{w}) \approx f(\mathbf{w}^t) + \partial f(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t)$$

but keeps ℓ as is to obtain the objective

$$\begin{aligned} \mathbf{w}^{t+1} &:= \arg \min_{\mathbf{w} \in \mathcal{W}} \ell(f(\mathbf{w}^t) + \partial f(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t)) \\ &= \arg \min_{\mathbf{w} \in \mathcal{W}} \ell(\partial f(\mathbf{w}^t)\mathbf{w} + f(\mathbf{w}^t) - \partial f(\mathbf{w}^t)\mathbf{w}^t) \\ &= \arg \min_{\mathbf{w} \in \mathcal{W}} \ell(J^t \mathbf{w} + \boldsymbol{\delta}^t), \end{aligned}$$

where we defined the shorthands $J^t := \partial f(\mathbf{w}^t)$ and $\boldsymbol{\delta}^t := f(\mathbf{w}^t) - \partial f(\mathbf{w}^t)\mathbf{w}^t$. We call $\ell(J^t \mathbf{w} + \boldsymbol{\delta}^t)$ the **partial linearization** of $L = \ell \circ f$ at \mathbf{w}^t , as opposed to the full linearization of L used in gradient descent.

Since the composition of a convex function and of linear function is convex, this objective is **convex** even if $L(\mathbf{w})$ is nonconvex. In practice, we often add a proximity term as regularization to define

$$\mathbf{w}^{t+1} := \arg \min_{\mathbf{w} \in \mathcal{W}} \ell(J^t \mathbf{w} + \boldsymbol{\delta}^t) + \frac{\eta^t}{2} \|\mathbf{w} - \mathbf{w}^t\|_2^2. \quad (17.3)$$

We can see this update as an approximation of the **proximal point** update

$$\arg \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}) + \frac{\eta^t}{2} \|\mathbf{w} - \mathbf{w}^t\|_2^2,$$

where $L(\mathbf{w})$ has been replaced by its partial linearization. Solving Eq. (17.3) using gradient-based solvers requires to compute the gradient

of $\mathbf{w} \mapsto \ell(J^t \mathbf{w} + \delta^t)$, which is $\mathbf{w} \mapsto (J^t)^* \nabla \ell(J^t \mathbf{w} + \delta^t)$. Computing this gradient by autodiff therefore requires to perform a forward pass to compute the JVP $J^t \mathbf{w}$ and a backward pass to compute the VJP $(J^t)^* \nabla \ell(\mathbf{z})$. See Section 2.3 for an introduction to these operators and Chapter 8 for an introduction to autodiff.

The Gauss-Newton method with arbitrary convex outer loss is often called modified Gauss-Newton (Nesterov, 2007) or prox-linear (Drusvyatskiy and Paquette, 2019). The classical Gauss-Newton and Levenberg-Marquardt (Levenberg, 1944; Marquardt, 1963) methods originate from nonlinear least-squares and are recovered when $\ell(\mathbf{z})$ is quadratic (Kelley, 1995), such as $\ell(\mathbf{z}) := \frac{1}{2} \|\mathbf{z} - \mathbf{y}\|_2^2$, for \mathbf{y} some reference target. The Gauss-Newton method corresponds classically to not using regularization (i.e., $\eta^t = 0$) and the Levenberg-Marquardt method uses regularization (usually called damping, potentially changing η^t across iterations). See e.g., (Messerer *et al.*, 2021), for a survey of different variants.

17.2.2 With approximate outer function

Another variant of the Gauss-Newton method consists in replacing the convex loss ℓ with its quadratic approximation around $\mathbf{z}^t := f(\mathbf{w}^t)$,

$$q^t(\mathbf{z}) := \ell(\mathbf{z}^t) + \langle \nabla \ell(\mathbf{z}^t), \mathbf{z} - \mathbf{z}^t \rangle + \frac{1}{2} \langle \mathbf{z} - \mathbf{z}^t, \nabla^2 \ell(\mathbf{z}^t)(\mathbf{z} - \mathbf{z}^t) \rangle \approx \ell(\mathbf{z})$$

to define the update

$$\mathbf{w}^{t+1} := \arg \min_{\mathbf{w} \in \mathcal{W}} q^t(J^t \mathbf{w} + \delta^t) + \frac{\eta^t}{2} \|\mathbf{w} - \mathbf{w}^t\|_2^2.$$

Notice that ℓ has been replaced by its quadratic approximation q^t . This objective is always a **convex quadratic**, unlike the objective of the Newton method in Eq. (17.1), which is a priori a **nonconvex quadratic**, if f is nonlinear. Simple calculations show that

$$\begin{aligned} \mathbf{w}^{t+1} = \arg \min_{\mathbf{w} \in \mathcal{W}} & L(\mathbf{w}^t) + \langle \mathbf{q}^t, J^t(\mathbf{w} - \mathbf{w}^t) \rangle \\ & + \frac{1}{2} \langle \mathbf{w} - \mathbf{w}^t, (J^t)^* Q^t J^t(\mathbf{w} - \mathbf{w}^t) \rangle + \frac{\eta^t}{2} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \end{aligned}$$

where $\mathbf{q}^t := \nabla \ell(f(\mathbf{w}^t)) \in \mathcal{M} = \mathbb{R}^Z$, $Q^t := \nabla^2 \ell(f(\mathbf{w}^t)) \in \mathbb{R}^{Z \times Z}$. The closed form solution is

$$\begin{aligned} \mathbf{w}^{t+1} &= \mathbf{w}^t - ((J^t)^* Q^t J^t + \eta^t \mathbf{I})^{-1} (J^t)^* \mathbf{q}^t \\ &= \mathbf{w}^t - (\nabla_{\text{GN}}^2(\ell \circ f)(\mathbf{w}^t) + \eta^t \mathbf{I})^{-1} \nabla L(\mathbf{w}^t), \end{aligned}$$

where we used the (generalized) **Gauss-Newton matrix** of $L = \ell \circ f$, defined in Section 9.2.

17.2.3 Linesearch

Similarly to Newton's method, the iterates of a Gauss-Newton method may diverge when used alone. However, the direction $-(\nabla_{\text{GN}}^2(\ell \circ f)(\mathbf{w}^t) + \eta^t \mathbf{I})^{-1} \nabla L(\mathbf{w}^t)$ defines a descent direction for any $\eta^t > 0$ and can be combined with a stepsize γ^t (typically chosen using a linesearch) to obtain iterates of the form

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \gamma^t (\nabla_{\text{GN}}^2(\ell \circ f)(\mathbf{w}^t) + \eta^t \mathbf{I})^{-1} \nabla L(\mathbf{w}^t).$$

17.2.4 Stochastic Gauss-Newton

In deep learning, the objective generally consists in an expectation over samples of the composition between a loss function and a network function:

$$L(\mathbf{w}) = \mathbb{E}_{S \sim \rho} [L(\mathbf{w}; S)] = \mathbb{E}_{(X, Y) \sim \rho} [\ell(f(\mathbf{w}; X); Y)]$$

where $S = (X, Y)$ denotes a sample pair of input X with associated label Y . In that case, as already studied in Section 9.2, the Gauss-Newton matrix $\nabla_{\text{GN}}^2 L$ is the expectation of the individual Gauss-Newton matrices

$$\begin{aligned} \nabla_{\text{GN}}^2 L(\mathbf{w}; \mathbf{x}, \mathbf{y}) &:= \partial f(\mathbf{w}; \mathbf{x})^\top \nabla^2 \ell(f(\mathbf{w}; \mathbf{x})) \partial f(\mathbf{w}; \mathbf{x}), \\ \nabla_{\text{GN}}^2 L(\mathbf{w}) &:= \mathbb{E}_{(X, Y) \sim \rho} [\nabla_{\text{GN}}^2 L(\mathbf{w}; \mathbf{x}, \mathbf{y})]. \end{aligned}$$

We can estimate the gradient and the Gauss-Newton matrix by, respectively, $\mathbf{g}(\mathbf{w}; S) \approx \nabla L(\mathbf{w})$, and $G(\mathbf{w}; S') \approx \nabla_{\text{GN}}^2 L(\mathbf{w})$ for $S, S' \sim \rho$ or using mini-batch approximations. A **stochastic** Gauss-Newton method therefore performs iterates

$$\mathbf{w}^{t+1} := \mathbf{w}^t - \gamma^t (G(\mathbf{w}; S') + \eta^t \mathbf{I})^{-1} \mathbf{g}(\mathbf{w}, S),$$

for $\eta^t \geq 0$ and γ^t fixed or selected to satisfy some criterion.

17.3 Natural gradient descent

Natural gradient descent (Amari, 1998) follows a similar principle as gradient descent: linearize the objective around the current iterate and minimize this approximation together with a proximity term. It differs from gradient descent in the choice of the proximity term: rather than using a squared Euclidean distance between the **parameters**, it uses a Kullback-Leibler divergence between the **probability distributions** these parameters define.

Negative log-likelihood

We consider objectives of the form

$$\min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}) = \mathbb{E}_{S \sim \rho} [L(\mathbf{w}; S)] = \mathbb{E}_{S \sim \rho} [-\log q_{\mathbf{w}}(S)],$$

where ρ is an unknown data distribution (but from which we can sample) and where $q_{\mathbf{w}}$ is a probability distribution parameterized by \mathbf{w} . As reviewed in Chapter 3, the negative log-likelihood can be used as a loss function (many loss functions can be seen from this perspective, including the squared and logistic loss functions). In the unsupervised setting, where $S = Y$, we simply use $q_{\mathbf{w}}(Y)$ as is. In the supervised setting, where $S = (X, Y)$, we use the product rule $\mathbb{P}(X, Y) = \mathbb{P}(X)\mathbb{P}(Y|X)$ to parameterize $q_{\mathbf{w}}(S)$ as

$$q_{\mathbf{w}}(\mathbf{x}, \mathbf{y}) := \rho_X(\mathbf{x})p_{\boldsymbol{\theta}}(\mathbf{y}),$$

where ρ_X is the marginal distribution for X , $p_{\boldsymbol{\theta}}(\mathbf{y})$ is the PMF/PDF of a probability distribution and $\boldsymbol{\theta} = f(\mathbf{w}; \mathbf{x})$ is for instance a neural network with parameters $\mathbf{w} \in \mathcal{W}$ and input $\mathbf{x} \in \mathcal{X}$.

17.3.1 Variational perspective

Natural gradient descent is motivated by updates of the form

$$\mathbf{w}^{t+1} = \arg \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}^t) + \langle \nabla L(\mathbf{w}^t), \mathbf{w} - \mathbf{w}^t \rangle + \text{KL}(q_{\mathbf{w}^t}, q_{\mathbf{w}}),$$

where $\text{KL}(p, q) := \int p(\mathbf{z}) \log \frac{p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z}$ is the Kullback-Leibler (KL) divergence. Unlike gradient descent, the proximity term is therefore between

the current distribution $q_{\mathbf{w}^t}$ and a candidate probability distribution $q_{\mathbf{w}}$. The above problem is intractable in general, as the KL may not have a closed form. Nevertheless, its quadratic approximation can be shown (Amari, 1998) to admit a simple form,

$$\text{KL}(q_{\mathbf{w}^t}, q_{\mathbf{w}}) \approx \frac{1}{2} \langle \mathbf{w} - \mathbf{w}^t, \nabla_{\mathbb{F}}^2 L(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) \rangle$$

where we used the Fisher information matrix $\nabla_{\mathbb{F}}^2 L(\mathbf{w})$, studied in Section 9.3. Equipped with this quadratic approximation of the KL divergence, natural gradient descent amounts to compute iterates as

$$\begin{aligned} \mathbf{w}^{t+1} := \arg \min_{\mathbf{w} \in \mathcal{W}} & L(\mathbf{w}^t) + \langle \nabla L(\mathbf{w}^t), \mathbf{w} - \mathbf{w}^t \rangle \\ & + \frac{1}{2} \langle \mathbf{w} - \mathbf{w}^t, \nabla_{\mathbb{F}}^2 L(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) \rangle + \frac{\eta^t}{2} \|\mathbf{w} - \mathbf{w}^t\|_2^2, \end{aligned}$$

where a quadratic proximity-term was added to ensure a unique solution. This is a strictly convex problem as $\nabla_{\mathbb{F}}^2 L(\mathbf{w}^t)$ is positive semi-definite. The closed-form solution is

$$\mathbf{w}^{t+1} = \mathbf{w}^t - (\nabla_{\mathbb{F}}^2 L(\mathbf{w}^t) + \eta^t \mathbf{I})^{-1} \nabla L(\mathbf{w}^t).$$

Because the Gauss-Newton and Fisher information matrices are equivalent when p_{θ} is an exponential family distribution (Proposition 9.6), the Gauss-Newton and natural gradient methods coincide in this case.

17.3.2 Stochastic natural gradient descent

In practice, we may not have access to $\nabla L(\mathbf{w}^t)$ in closed form as it is an expectation over ρ . Moreover, $\nabla_{\mathbb{F}}^2 L(\mathbf{w}^t)$ may not be computable in closed form either. To estimate the Fisher information matrix, we can use that (see Section 9.3) using the shorthand $\boldsymbol{\theta} := f(\mathbf{w}, X)$,

$$\nabla_{\mathbb{F}}^2 L(\mathbf{w}) = \mathbb{E}_{X \sim \rho_X} \mathbb{E}_{Y \sim p_{\theta}} [\nabla L(\mathbf{w}; X, Y) \otimes \nabla L(\mathbf{w}; X, Y)].$$

We can then build estimates $\mathbf{g}(\mathbf{w}^t; S) \approx \nabla L(\mathbf{w}^t, S)$ and $F(\mathbf{w}^t; S') \approx \nabla_{\mathbb{F}}^2 L(\mathbf{w}^t)$ for S sampled from ρ and S' sampled from $q_{\mathbf{w}^t}(\mathbf{x}, \mathbf{y}) = p_X(\mathbf{x})\rho_{\theta}(\mathbf{y})$. A stochastic natural gradient descent can then be implemented as

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \gamma^t (F(\mathbf{w}^t; S') + \eta^t \mathbf{I})^{-1} \mathbf{g}(\mathbf{w}^t; S),$$

where γ^t is a stepsize, possibly chosen by linesearch.

In deep learning, the product with the inverse Fisher or Gauss-Newton matrices can remain costly to compute. Several approximations have been proposed, such as KFAC (Martens and Grosse, 2015; Botev *et al.*, 2017), which uses a computationally efficient structural approximation to these matrices.

17.4 Quasi-Newton methods

17.4.1 BFGS

A celebrated example of **quasi-Newton** method is the BFGS method (Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970), whose acronym follows from its author names. The rationale of the BFGS update stems once again from a variational viewpoint. We wish to build a simple quadratic model of the objective $h^t(\mathbf{w}) = L(\mathbf{w}^t) + \langle \nabla L(\mathbf{w}^t), \mathbf{w} - \mathbf{w}^t \rangle + \frac{1}{2} \langle \mathbf{w} - \mathbf{w}^t, Q^t(\mathbf{w} - \mathbf{w}^t) \rangle$ for some Q^t built along the iterations rather than taken as $\nabla^2 L(\mathbf{w}^t)$. One desirable property of such quadratic model would be that its gradients at consecutive iterates match the gradients of the original function, i.e., $\nabla h^t(\mathbf{w}^t) = \nabla L(\mathbf{w}^t)$ and $\nabla h^t(\mathbf{w}^{t-1}) = \nabla L(\mathbf{w}^{t-1})$. A simpler condition, called the **secant condition** consists in considering the differences of these vectors, that is, ensuring that

$$\begin{aligned} \nabla h^t(\mathbf{w}^t) - \nabla h^t(\mathbf{w}^{t-1}) &= \nabla L(\mathbf{w}^t) - \nabla L(\mathbf{w}^{t-1}) \\ \iff Q^t(\mathbf{w}^t - \mathbf{w}^{t-1}) &= \nabla L(\mathbf{w}^t) - \nabla L(\mathbf{w}^{t-1}) \\ \iff \mathbf{w}^t - \mathbf{w}^{t-1} &= B^t(\nabla L(\mathbf{w}^t) - \nabla L(\mathbf{w}^{t-1})), \end{aligned}$$

for $B^t = (Q^t)^{-1}$. Building B^t , a surrogate of the inverse of the Hessian satisfying the secant equation, can then be done as

$$B^{t+1} := \left(\mathbf{I} - \rho^t \mathbf{s}^t (\mathbf{y}^t)^\top \right) B^t \left(\mathbf{I} - \rho^t \mathbf{s}^t (\mathbf{y}^t)^\top \right) + \rho^t \mathbf{s}^t (\mathbf{s}^t)^\top$$

where

$$\begin{aligned} \mathbf{s}^t &:= \mathbf{w}^{t+1} - \mathbf{w}^t \\ \mathbf{y}^t &:= \nabla L(\mathbf{w}^{t+1}) - \nabla L(\mathbf{w}^t) \\ \rho^t &:= \frac{1}{\langle \mathbf{s}^t, \mathbf{y}^t \rangle}. \end{aligned}$$

A typical implementation of BFGS stores $B_t \in \mathbb{R}^{P \times P}$ in memory, which is prohibitive when P is large.

17.4.2 Limited-memory BFGS

In practice, the limited-memory counterpart of BFGS, called LBFGS (Liu and Nocedal, 1989), is often preferred. The key observation of LBFGS is that we do not need to materialize B_t in memory: we only need to multiply it with the gradient $\nabla L(\mathbf{w}^t)$. That is, we can see B_t as a linear map. Fortunately, the product between B^t and any vector \mathbf{v} can be computed efficiently if we store $(\mathbf{s}^1, \mathbf{y}^1, \rho^1), \dots, (\mathbf{s}^t, \mathbf{y}^t, \rho^t)$ in memory. In practice, a small history of past values is used to reduce memory and computational cost. Because LBFGS has the benefits of second-order-like methods with much reduced cost, it has become a de-facto algorithm, outperforming most other algorithms for medium-scale problems without particular structure (Liu and Nocedal, 1989).

17.5 Approximate Hessian diagonal inverse preconditioners

One application of the approximations of the Hessian diagonal developed in Section 9.7 is to obtain cheap approximations of the Hessian diagonal inverse,

$$B^t := \mathbf{diag}(|H_{11}^t|^{-1}, \dots, |H_{PP}^t|^{-1}).$$

Such a scaling would for instance be sufficient to make the quadratic example presented in Fig. 17.1 work. Many optimization algorithms, including the popular ADAM, can be viewed as using a preconditioner that approximates the inverse of the Hessian's diagonal.

17.6 Summary

- We reviewed Newton's method, the Gauss-Newton method, natural gradient descent, quasi-Newton methods and preconditioning methods.
- We adopted a variational viewpoint, where the method's next iterate is computed as the solution of a trade-off between mini-

mizing an approximation of the function (linear, partially linear, quadratic) and a proximity term (squared Euclidean, KL).

- All methods were shown to use iterates of the form

$$\mathbf{w}^{t+1} := \mathbf{w}^t - \gamma^t B^t \nabla L(\mathbf{w}^t)$$

but have different trade-offs between the cost it takes to evaluate $B^t \nabla L(\mathbf{w}^t)$ and the richness of the information used about L .

18

Duality

In this chapter, we review duality principles in optimization.

18.1 Dual norms

We introduce in this section dual norms, since they are useful in this book.

Definition 18.1 (Dual norms). Given a norm $\|\mathbf{u}\|$, its dual is

$$\|\mathbf{v}\|_* := \max_{\|\mathbf{u}\| \leq 1} \langle \mathbf{u}, \mathbf{v} \rangle.$$

Therefore, the dual norm of $\|\cdot\|$ is the **support function** of the unit ball induced by the norm $\|\cdot\|$,

$$B_{\|\cdot\|} := \{\mathbf{u} \in \mathbb{R}^D : \|\mathbf{u}\| \leq 1\}.$$

We give examples of pairs of dual norms below.

Example 18.1 (Dual norm of p -norms). The p -norm is defined by

$$\|\mathbf{u}\|_p := \left(\sum_{j=1}^D |u_j|^p \right)^{1/p}.$$

Its dual is $\|\mathbf{v}\|_q$ where q is such that $\frac{1}{p} + \frac{1}{q} = 1$. For instance, the dual norm of the 2-norm is itself, since $\frac{1}{2} + \frac{1}{2} = 1$. The 1-norm and the ∞ -norm are dual of each other, since $\frac{1}{1} + \frac{1}{\infty} = 1$.

The definition of dual norm implies a generalization of **Cauchy–Schwarz’s inequality**: for all $\mathbf{u}, \mathbf{v} \in \mathbb{R}^D$

$$|\langle \mathbf{u}, \mathbf{v} \rangle| \leq \|\mathbf{u}\|_* \|\mathbf{v}\|.$$

See, e.g., Beck (2017, Lemma 1.4).

Proposition 18.1 (Conjugate of norms and squared norms). We know that the conjugate of the support function is the indicator function. Therefore, if $f(\mathbf{u}) = \|\mathbf{u}\|$, then

$$f^*(\mathbf{v}) = \iota_{B_{\|\cdot\|}}(\mathbf{v}) = \begin{cases} 0 & \text{if } \|\mathbf{v}\|_* \leq 1 \\ \infty & \text{otherwise} \end{cases}.$$

On the other hand, if $f(\mathbf{u}) = \frac{1}{2}\|\mathbf{u}\|^2$, then

$$f^*(\mathbf{v}) = \frac{1}{2}\|\mathbf{v}\|_*^2.$$

18.2 Fenchel duality

We consider in this section standard objectives of the form

$$\min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}) := \min_{\mathbf{w} \in \mathcal{W}} \ell(f(\mathbf{w})) + R(\mathbf{w}),$$

where $f: \mathcal{W} \rightarrow \mathcal{M}$, $\ell: \mathcal{M} \rightarrow \mathbb{R}$ and $R: \mathcal{W} \rightarrow \mathbb{R}$. We first show that the minimization of this objective, called the **primal**, can be lower bounded by a **concave** maximization objective, called the **dual**, even if the primal is nonconvex.

Proposition 18.2 (Weak duality). Let $f: \mathcal{W} \rightarrow \mathcal{M}$ (potentially non-linear), $\ell: \mathcal{M} \rightarrow \mathbb{R}$ (potentially nonconvex) and $R: \mathcal{W} \rightarrow \mathbb{R}$ (potentially nonconvex). Then

$$\min_{\mathbf{w} \in \mathcal{W}} \ell(f(\mathbf{w})) + R(\mathbf{w}) \geq \max_{\boldsymbol{\alpha} \in \mathcal{M}} -R^f(\boldsymbol{\alpha}) - \ell^*(-\boldsymbol{\alpha}),$$

where we used the conjugate

$$\ell^*(-\boldsymbol{\alpha}) := \max_{\boldsymbol{\theta} \in \mathcal{M}} \langle -\boldsymbol{\alpha}, \boldsymbol{\theta} \rangle - \ell(\boldsymbol{\theta})$$

and the “generalized conjugate”

$$R^f(\boldsymbol{\alpha}) := \max_{\mathbf{w} \in \mathcal{W}} \langle \boldsymbol{\alpha}, f(\mathbf{w}) \rangle - R(\mathbf{w}).$$

Moreover, ℓ^* and R^f are both convex functions.

We emphasize that the result in Proposition 18.2 is fully general, in the sense that it does not assume the linearity of f or the convexity of ℓ and R . The caveat, of course, is that R^f and ℓ^* are difficult to compute in general, if f is nonlinear, and if ℓ and R are nonconvex.

Proof.

$$\begin{aligned} & \min_{\mathbf{w} \in \mathcal{W}} \ell(f(\mathbf{w})) + R(\mathbf{w}) \\ &= \min_{\substack{\mathbf{w} \in \mathcal{W} \\ \boldsymbol{\theta} \in \mathcal{M}}} \ell(\boldsymbol{\theta}) + R(\mathbf{w}) \quad \text{s.t.} \quad \boldsymbol{\theta} = f(\mathbf{w}) \\ &= \min_{\substack{\mathbf{w} \in \mathcal{W} \\ \boldsymbol{\theta} \in \mathcal{M}}} \max_{\boldsymbol{\alpha} \in \mathcal{M}} \ell(\boldsymbol{\theta}) + R(\mathbf{w}) + \langle \boldsymbol{\alpha}, \boldsymbol{\theta} - f(\mathbf{w}) \rangle \\ &\geq \max_{\boldsymbol{\alpha} \in \mathcal{M}} \min_{\substack{\mathbf{w} \in \mathcal{W} \\ \boldsymbol{\theta} \in \mathcal{M}}} \ell(\boldsymbol{\theta}) + R(\mathbf{w}) + \langle \boldsymbol{\alpha}, \boldsymbol{\theta} - f(\mathbf{w}) \rangle \\ &= \max_{\boldsymbol{\alpha} \in \mathcal{M}} \min_{\mathbf{w} \in \mathcal{W}} \langle \boldsymbol{\alpha}, -f(\mathbf{w}) \rangle + R(\mathbf{w}) + \min_{\boldsymbol{\theta} \in \mathcal{M}} \ell(\boldsymbol{\theta}) + \langle \boldsymbol{\alpha}, \boldsymbol{\theta} \rangle \\ &= \max_{\boldsymbol{\alpha} \in \mathcal{M}} -\max_{\mathbf{w} \in \mathcal{W}} \langle \boldsymbol{\alpha}, f(\mathbf{w}) \rangle - R(\mathbf{w}) - \max_{\boldsymbol{\theta} \in \mathcal{M}} \langle -\boldsymbol{\alpha}, \boldsymbol{\theta} \rangle - \ell(\boldsymbol{\theta}) \\ &= \max_{\boldsymbol{\alpha} \in \mathcal{M}} -R^f(\boldsymbol{\alpha}) - \ell^*(-\boldsymbol{\alpha}). \end{aligned}$$

□

In the case when $f(\mathbf{w}) = A\mathbf{w}$, where A is a linear map, and when both ℓ and R are convex, we can state a much stronger result.

Proposition 18.3 (Strong duality). Let A be a linear map from \mathcal{W} to \mathcal{M} . Let $\ell: \mathcal{M} \rightarrow \mathbb{R}$ and $R: \mathcal{W} \rightarrow \mathbb{R}$ be convex functions. Let A^* denote the adjoint of A (Section 2.3). Then,

$$\min_{\mathbf{w} \in \mathcal{W}} \ell(A\mathbf{w}) + R(\mathbf{w}) = \max_{\boldsymbol{\alpha} \in \mathcal{M}} -R^*(A^*\boldsymbol{\alpha}) - \ell^*(-\boldsymbol{\alpha}).$$

Furthermore, the primal solution satisfies

$$\mathbf{w}^* \in \arg \max_{\mathbf{w} \in \mathcal{W}} \langle A\boldsymbol{\alpha}^*, \mathbf{w} \rangle - R(\mathbf{w}).$$

When R is strictly convex, the primal solution is uniquely determined by

$$\mathbf{w}^* = \nabla R^*(A^*\boldsymbol{\alpha}^*).$$

Proof. Since $f(\mathbf{w}) = A\mathbf{w}$, we have

$$\begin{aligned} R^f(\boldsymbol{\alpha}) &:= \max_{\mathbf{w} \in \mathcal{W}} \langle \boldsymbol{\alpha}, f(\mathbf{w}) \rangle - R(\mathbf{w}) \\ &= \max_{\mathbf{w} \in \mathcal{W}} \langle \boldsymbol{\alpha}, A\mathbf{w} \rangle - R(\mathbf{w}) \\ &= \max_{\mathbf{w} \in \mathcal{W}} \langle A^*\boldsymbol{\alpha}, \mathbf{w} \rangle - R(\mathbf{w}) \\ &= R^*(A^*\boldsymbol{\alpha}). \end{aligned}$$

Furthermore, the inequality in the proof of Proposition 18.2 is an equality, since the min max is that of a convex-concave function. \square

The maximization problem in Proposition 18.3 is called the **Fenchel dual**. By strong duality, the value of the maximum and the value of the minimum are equal. We can therefore choose to equivalently solve the dual instead of the primal. This can be advantageous when the space \mathcal{M} is smaller than \mathcal{W} .

We now apply the Fenchel dual to obtain the dual of regularized multiclass linear classification.

Table 18.1: Examples of loss conjugates. For regression losses (squared, absolute), where $\mathbf{y}_i \in \mathbb{R}^M$, we define $\mathbf{t}_i = \phi(\mathbf{y}_i) = \mathbf{y}_i$. For classification losses (logistic, perceptron, hinge), where $y_i \in [M]$, we define $\mathbf{t}_i = \phi(y_i) = \mathbf{e}_{y_i}$. To simplify some expressions, we defined the change of variable $\boldsymbol{\mu}_i := \mathbf{y}_i - \boldsymbol{\alpha}_i$.

	$\ell_i(\boldsymbol{\theta}_i)$	$\ell_i^*(-\boldsymbol{\alpha}_i)$
Squared	$\frac{1}{2}\ \boldsymbol{\theta}_i - \mathbf{t}_i\ _2^2$	$\frac{1}{2}\ \boldsymbol{\alpha}_i\ _2^2 - \langle \mathbf{t}_i, \boldsymbol{\alpha}_i \rangle$
Absolute	$\ \boldsymbol{\theta}_i - \mathbf{t}_i\ _1$	$\iota_{[-1,1]^M}(\boldsymbol{\alpha}_i) - \langle \mathbf{t}_i, \boldsymbol{\alpha}_i \rangle$
Logistic	$\text{LSE}(\boldsymbol{\theta}_i) - \langle \boldsymbol{\theta}, \mathbf{t}_i \rangle$	$\langle \boldsymbol{\mu}_i, \log \boldsymbol{\mu}_i \rangle + \iota_{\Delta^M}(\boldsymbol{\mu}_i)$
Perceptron	$\max_{i \in [M]} \theta_i - \theta_y$	$\iota_{\Delta^M}(\boldsymbol{\mu}_i)$
Hinge	$\max_{i \in [M]} [i \neq y] + \theta_i - \theta_y$	$\iota_{\Delta^M}(\boldsymbol{\mu}_i) - \langle \mathbf{1} - \mathbf{t}_i, \boldsymbol{\mu}_i \rangle$

Example 18.2 (Sum of separable loss functions). When the loss is $\ell(\boldsymbol{\theta}) := \sum_{i=1}^N \ell_i(\boldsymbol{\theta}_i)$, where $\boldsymbol{\theta} = A\mathbf{w} = (A_1\mathbf{w}, \dots, A_N\mathbf{w}) \in \mathcal{M}^N$ and A_i is a linear map from \mathcal{W} to \mathcal{M} , we obtain

$$\min_{\mathbf{w} \in \mathcal{W}} \sum_{i=1}^N \ell_i(A_i\mathbf{w}) + R(\mathbf{w}) = \max_{\boldsymbol{\alpha} \in \mathcal{M}^N} -R(A^*\boldsymbol{\alpha}) - \sum_{i=1}^N \ell_i^*(-\boldsymbol{\alpha}_i),$$

where $A^*\boldsymbol{\alpha} = (A_1^*\boldsymbol{\alpha}_1, \dots, A_N^*\boldsymbol{\alpha}_N)$. Typically, we define

$$A_i\mathbf{w} := \mathbf{W}\mathbf{x}_i,$$

where $\mathbf{W} \in \mathbb{R}^{M \times D}$ is a reshaped version of $\mathbf{w} \in \mathcal{W}$, $\mathbf{x}_i \in \mathbb{R}^D$ is a training sample, and M is the number of classes. In this case, we then have

$$A_i^*\boldsymbol{\alpha}_i = \boldsymbol{\alpha}_i \mathbf{x}_i^\top.$$

Examples of loss function conjugates are given in Table 18.1.

18.3 Bregman divergences

Bregman divergences are a measure of difference between two points.

Definition 18.2 (Bregman divergence). The Bregman divergence gen-

erated by a differentiable convex function $f: \mathbb{R}^D \rightarrow \mathbb{R}$ is

$$\begin{aligned} B_f(\mathbf{u}, \mathbf{v}) &:= f(\mathbf{u}) - f(\mathbf{v}) - \langle \nabla f(\mathbf{v}), \mathbf{u} - \mathbf{v} \rangle \\ &= \langle \nabla f(\mathbf{v}), \mathbf{v} \rangle - f(\mathbf{v}) - [\langle \nabla f(\mathbf{v}), \mathbf{u} \rangle - f(\mathbf{u})], \end{aligned}$$

where $\mathbf{u}, \mathbf{v} \in \text{dom}(f)$.

Intuitively, the Bregman divergence is the difference between $f(\mathbf{u})$ and its linearization $\mathbf{u} \mapsto f(\mathbf{v}) + \langle \nabla f(\mathbf{v}), \mathbf{u} - \mathbf{v} \rangle$ around \mathbf{v} . This is illustrated in Fig. 18.1.

Example 18.3 (Examples of Bregman divergences). If $f(\mathbf{u}) = \frac{1}{2}\|\mathbf{u}\|_2^2$, where $\text{dom}(f) = \mathbb{R}^D$, then

$$B_f(\mathbf{u}, \mathbf{v}) = \frac{1}{2}\|\mathbf{u} - \mathbf{v}\|_2^2,$$

the **squared Euclidean distance**. If $f(\mathbf{u}) = \langle \mathbf{u}, \log \mathbf{u} \rangle$, where $\text{dom}(f) = \mathbb{R}_+^D$, then

$$B_f(\mathbf{u}, \mathbf{v}) = \sum_{j=1}^D u_j \log \frac{u_j}{v_j} - \sum_{j=1}^D u_j + \sum_{j=1}^D v_j,$$

the (generalized) **Kullback-Leibler divergence**.

Properties

Bregman divergences enjoy several useful properties.

Proposition 18.4 (Properties of Bregman divergences). Let $f: \mathbb{R}^D \rightarrow \mathbb{R}$ be a differentiable convex function.

1. **Non-negativity:** $B_f(\mathbf{u}, \mathbf{v}) \geq 0$ for all $\mathbf{u}, \mathbf{v} \in \text{dom}(f)$.
2. **Positivity:** $B_f(\mathbf{u}, \mathbf{v}) = 0$ if and only if $\mathbf{u} = \mathbf{v}$ (when f is strictly convex).
3. **Convexity:** $B_f(\mathbf{u}, \mathbf{v})$ is convex in \mathbf{u} .
4. **Dual-space form:** $B_f(\mathbf{u}, \mathbf{v}) = B_{f^*}(\mathbf{b}, \mathbf{a})$, where $\mathbf{b} = \nabla f(\mathbf{v}) \in$

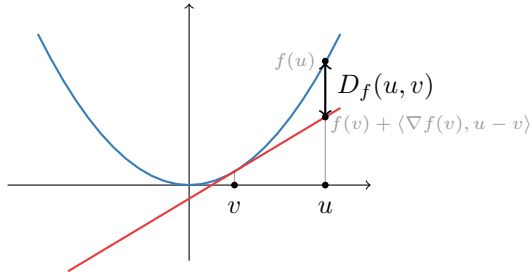


Figure 18.1: The Bregman divergence generated by f is the difference between $f(\mathbf{u})$ and its linearization around \mathbf{v} .

dom(f^*) and $\mathbf{a} = \nabla f(\mathbf{u}) \in \text{dom}(f^*)$.

Proof. The properties follow immediately from the convexity of $f(\mathbf{u})$.

1. From Definition 15.6.
2. From the unicity of minimizers.
3. From the fact that $\mathbf{u} \mapsto B_f(\mathbf{u}, \mathbf{v})$ is the sum of $f(\mathbf{u})$ and a linear function of \mathbf{u} .

□

The Bregman divergence can be used to define natural generalizations of the Euclidean projection and proximal operators, reviewed in Section 16.3 and Section 16.4.

Definition 18.3 (Bregman proximal and projection operators). Let $\mathbf{v} \in \text{dom}(f)$. The Bregman proximal operator is

$$\text{bprox}_{f,g}(\mathbf{v}) := \arg \min_{\mathbf{u} \in \text{dom}(f) \cap \text{dom}(g)} B_f(\mathbf{u}, \mathbf{v}) + g(\mathbf{u}).$$

In particular, the Bregman projection onto $\mathcal{C} \subseteq \text{dom}(f)$ is

$$\text{bproj}_{f,\mathcal{C}}(\mathbf{v}) := \arg \min_{\mathbf{u} \in \mathcal{C}} B_f(\mathbf{u}, \mathbf{v}).$$

It turns out that these operators are intimately connected to the gradient mapping of the convex conjugate.

Proposition 18.5 (Link with conjugate's gradient). If $\Omega = f + g$, then for all $\boldsymbol{\theta} \in \text{dom}(f^*)$

$$\nabla \Omega^*(\boldsymbol{\theta}) = \text{bprox}_{f,g}(\nabla f^*(\boldsymbol{\theta})).$$

In particular, if $\Omega = f + \iota_C$, then for all $\boldsymbol{\theta} \in \text{dom}(f^*)$

$$\nabla \Omega^*(\boldsymbol{\theta}) = \text{bproj}_{f,C}(\nabla f^*(\boldsymbol{\theta})).$$

We give two examples below.

Example 18.4 (Bregman projections on the simplex). If $f(\mathbf{u}) = \frac{1}{2}\|\mathbf{u}\|_2^2$, then

$$\text{bproj}_{f,\Delta^D}(\mathbf{v}) = \arg \min_{\mathbf{u} \in \Delta^D} \frac{1}{2}\|\mathbf{u} - \mathbf{v}\|_2^2.$$

If $f(\mathbf{u}) = \langle \mathbf{u}, \log \mathbf{u} - \mathbf{1} \rangle$, then

$$\text{bproj}_{f,\Delta^D}(\mathbf{v}) = \arg \min_{\mathbf{u} \in \mathbb{R}_+^D} \text{KL}(\mathbf{u}, \mathbf{v}) = \text{softmax}(\boldsymbol{\theta}),$$

where $\mathbf{v} = \nabla f^*(\boldsymbol{\theta}) = \exp(\boldsymbol{\theta})$.

Therefore, the softmax can be seen as a projection onto the probability simplex in the Kullback-Leibler divergence sense!

18.4 Fenchel-Young loss functions

We end this chapter with a brief review of the Fenchel-Young family of loss functions (Blondel *et al.*, 2020), which includes all loss functions in Table 18.1.

Definition 18.4 (Fenchel-Young loss). The Fenchel-Young loss function generated by Ω is

$$\ell_\Omega(\boldsymbol{\theta}, \mathbf{t}) := \Omega^*(\boldsymbol{\theta}) + \Omega(\mathbf{t}) - \langle \boldsymbol{\theta}, \mathbf{t} \rangle$$

where $\boldsymbol{\theta} \in \text{dom}(\Omega^*)$ and $\mathbf{t} \in \text{dom}(\Omega)$.

Typically, we set $\boldsymbol{\theta} = f(\mathbf{x}, \mathbf{w})$, where f is a model prediction function with parameters \mathbf{w} and $\mathbf{t} = \phi(\mathbf{y})$, where $\phi: \mathcal{Y} \rightarrow \text{dom}(\Omega)$. For instance,

suppose we work with categorical outputs $y \in [M]$. Then, we can set $\phi(y) = \mathbf{e}_y$, where \mathbf{e}_y is the one-hot encoding of y .

The important point to notice is that the Fenchel-Young loss is defined over arguments in **mixed spaces**: $\boldsymbol{\theta}$ belongs to the dual space, while \mathbf{t} belongs to the primal space. In fact, the Fenchel-Young loss is intimately connected to the Bregman divergence, since $B_\Omega(\mathbf{t}, \mathbf{v}) = \Omega^*(\boldsymbol{\theta}) + \Omega(\mathbf{t}) - \langle \boldsymbol{\theta}, \mathbf{t} \rangle$, if we set $\boldsymbol{\theta} = \nabla \Omega(\mathbf{v})$. The key properties of Fenchel-Young loss functions are summarized below.

Proposition 18.6 (Properties of Fenchel-Young loss functions).

1. **Non-negativity**: $\ell_\Omega(\boldsymbol{\theta}, \mathbf{t}) \geq 0$ for all $\boldsymbol{\theta} \in \text{dom}(\Omega^*)$ and $\mathbf{t} \in \text{dom}(\Omega)$.
2. **Positivity**: $\ell_\Omega(\boldsymbol{\theta}, \mathbf{t}) = 0$ if and only if $\nabla \Omega^*(\boldsymbol{\theta}) = \mathbf{t}$, assuming Ω is strictly convex.
3. **Convexity**: $\ell_\Omega(\boldsymbol{\theta}, \mathbf{t})$ is convex in $\boldsymbol{\theta}$ (regardless of Ω) and in \mathbf{t} (if Ω is convex)
4. **Relation with composite Bregman divergence**:

$$0 \leq \underbrace{B_\Omega(\mathbf{t}, \nabla \Omega^*(\boldsymbol{\theta}))}_{\text{possibly nonconvex in } \boldsymbol{\theta}} \leq \underbrace{\ell_\Omega(\boldsymbol{\theta}, \mathbf{t})}_{\text{convex in } \boldsymbol{\theta}}.$$

See Blondel *et al.* (2020) for an in-depth study of more properties.

18.5 Summary

- The convex conjugate serves as a powerful abstraction in **Fenchal duality**, decoupling the dual expression and function-specific terms.
- The convex conjugate is also tightly connected to **Bregman divergences** and can be used to derive the family of **Fenchel-Young loss** functions, which can be seen as primal-dual Bregman divergences.

References

- Abadi, M., A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, *et al.* (2016). “Tensorflow: Large-scale machine learning on heterogeneous distributed systems”. *arXiv preprint arXiv:1603.04467*.
- Abernethy, J., C. Lee, and A. Tewari. (2016). “Perturbation techniques in online learning and optimization”. *Perturbations, Optimization, and Statistics*. 233.
- Aji, S. M. and R. J. McEliece. (2000). “The generalized distributive law”. *IEEE transactions on Information Theory*. 46(2): 325–343.
- Amari, S.-I. (1998). “Natural gradient works efficiently in learning”. *Neural computation*. 10(2): 251–276.
- Ba, J. L., J. R. Kiros, and G. E. Hinton. (2016). “Layer normalization”. *arXiv preprint arXiv:1607.06450*.
- Bach, F., R. Jenatton, J. Mairal, G. Obozinski, *et al.* (2012). “Optimization with sparsity-inducing penalties”. *Foundations and Trends® in Machine Learning*. 4(1): 1–106.
- Ball, K., E. A. Carlen, and E. H. Lieb. (2002). “Sharp uniform convexity and smoothness inequalities for trace norms”. *Inequalities: Selecta of Elliott H. Lieb*: 171–190.
- Ball, W. W. R. (1960). *A short account of the history of mathematics*. Courier Corporation.

- Balog, M., N. Tripuraneni, Z. Ghahramani, and A. Weller. (2017). “Lost relatives of the Gumbel trick”. In: *International Conference on Machine Learning*. PMLR. 371–379.
- Barndorff-Nielsen, O. (2014). *Information and exponential families: in statistical theory*. John Wiley & Sons.
- Baston, R. A. and Y. Nakatsukasa. (2022). “Stochastic diagonal estimation: probabilistic bounds and an improved algorithm”. *arXiv preprint arXiv:2201.10684*.
- Baum, L. E. and T. Petrie. (1966). “Statistical inference for probabilistic functions of finite state Markov chains”. *The annals of mathematical statistics*. 37(6): 1554–1563.
- Baur, W. and V. Strassen. (1983). “The complexity of partial derivatives”. *Theoretical computer science*. 22(3): 317–330.
- Bauschke Heinz, H. and L. Combettes Patrick. (2017). *Convex Analysis and Monotone Operator Theory in Hilbert Spaces, 2011*. 2nd ed. 978–1.
- Baydin, A. G., B. A. Pearlmutter, A. A. Radul, and J. M. Siskind. (2018). “Automatic differentiation in machine learning: a survey”. *Journal of Machine Learning Research*. 18: 1–43.
- Baydin, A. G., B. A. Pearlmutter, D. Syme, F. Wood, and P. Torr. (2022). “Gradients without backpropagation”. *arXiv preprint arXiv:2202.08587*.
- Beck, A. (2017). *First-order methods in optimization*. SIAM.
- Beck, A. and M. Teboulle. (2012). “Smoothing and first order methods: A unified framework”. *SIAM Journal on Optimization*. 22(2): 557–580.
- Becker, S. and Y. Le Cun. (1988). “Improving the convergence of back-propagation learning with second order methods”. In: *Proceedings of the 1988 connectionist models summer school*. 29–37.
- Bekas, C., E. Kokiopoulou, and Y. Saad. (2007). “An estimator for the diagonal of a matrix”. *Applied numerical mathematics*. 57(11-12): 1214–1229.
- Bergstra, J., O. Breuleux, F. Bastien, P. Lamblin, R. Pascanu, G. Desjardins, J. Turian, D. Warde-Farley, and Y. Bengio. (2010). “Theano: a CPU and GPU math expression compiler”. In: *Proceedings of the Python for scientific computing conference (SciPy)*. Vol. 4. No. 3. Austin, TX. 1–7.

- Berthet, Q., M. Blondel, O. Teboul, M. Cuturi, J.-P. Vert, and F. Bach. (2020). “Learning with differentiable perturbed optimizers”. *Advances in neural information processing systems*. 33: 9508–9519.
- Blelloch, G. E. (1989). “Scans as primitive parallel operations”. *IEEE Transactions on computers*. 38(11): 1526–1538.
- Blondel, M. (2019). “Structured prediction with projection oracles”. *Advances in neural information processing systems*. 32.
- Blondel, M., Q. Berthet, M. Cuturi, R. Frostig, S. Hoyer, F. Llinares-López, F. Pedregosa, and J.-P. Vert. (2021). “Efficient and Modular Implicit Differentiation”. *arXiv preprint arXiv:2105.15183*.
- Blondel, M., A. F. Martins, and V. Niculae. (2020). “Learning with fenchel-young losses”. *The Journal of Machine Learning Research*. 21(1): 1314–1382.
- Bolte, J., R. Boustany, E. Pauwels, and B. Pesquet-Popescu. (2022). “On the complexity of nonsmooth automatic differentiation”. In: *The Eleventh International Conference on Learning Representations*.
- Bolte, J. and E. Pauwels. (2020). “A mathematical model for automatic differentiation in machine learning”. *Advances in Neural Information Processing Systems*. 33: 10809–10819.
- Botev, A., H. Ritter, and D. Barber. (2017). “Practical Gauss-Newton optimisation for deep learning”. In: *International Conference on Machine Learning*. 557–565.
- Boumal, N. (2023). *An introduction to optimization on smooth manifolds*. Cambridge University Press.
- Boyd, S. P. and L. Vandenberghe. (2004). *Convex optimization*. Cambridge university press.
- Bradbury, J., R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. (2018). *JAX: composable transformations of Python+NumPy programs*. Version 0.3.13. URL: <http://github.com/google/jax>.
- Braun, M. and M. Golubitsky. (1983). *Differential equations and their applications*. Vol. 2. Springer.

- Brockhoff, D., A. Auger, N. Hansen, D. V. Arnold, and T. Hohm. (2010). “Mirrored sampling and sequential selection for evolution strategies”. In: *Parallel Problem Solving from Nature, PPSN XI: 11th International Conference, Kraków, Poland, September 11-15, 2010, Proceedings, Part I* 11. Springer. 11–21.
- Broyden, C. G. (1970). “The convergence of a class of double-rank minimization algorithms 1. general considerations”. *IMA Journal of Applied Mathematics*. 6(1): 76–90.
- Brucker, P. (1984). “An $O(n)$ algorithm for quadratic knapsack problems”. *Operations Research Letters*. 3(3): 163–166.
- Butcher, J. C. (2016). *Numerical methods for ordinary differential equations*. John Wiley & Sons.
- Cajori, F. (1993). *A history of mathematical notations*. Vol. 1. Courier Corporation.
- Céa, J. (1986). “Conception optimale ou identification de formes, calcul rapide de la dérivée directionnelle de la fonction coût”. *M2AN-Modélisation mathématique et analyse numérique*. 20(3): 371–402.
- Chaudhuri, S. and A. Solar-Lezama. (2010). “Smooth interpretation”. *ACM Sigplan Notices*. 45(6): 279–291.
- Chen, R. T., Y. Rubanova, J. Bettencourt, and D. K. Duvenaud. (2018). “Neural ordinary differential equations”. *Advances in neural information processing systems*. 31.
- Chen, X., N. Kayal, A. Wigderson, *et al.* (2011). “Partial derivatives in arithmetic complexity and beyond”. *Foundations and Trends® in Theoretical Computer Science*. 6(1–2): 1–138.
- Clarke, F. H., Y. S. Ledyaev, R. J. Stern, and P. R. Wolenski. (2008). *Nonsmooth analysis and control theory*. Vol. 178. Springer Science & Business Media.
- Clarke, F. H. (1975). “Generalized gradients and applications”. *Transactions of the American Mathematical Society*. 205: 247–262.
- Cohn, D. L. (2013). *Measure theory*. Vol. 5. Springer.
- Condat, L. (2016). “Fast projection onto the simplex and the ℓ_1 ball”. *Mathematical Programming*. 158(1-2): 575–585.
- Dalrymple, D. (2016). *Differentiable Programming*. URL: <https://www.edge.org/response-detail/26794>.

- Dangel, F., F. Kunstner, and P. Hennig. (2019). “Backpack: Packing more into backprop”. *arXiv preprint arXiv:1912.10985*.
- Davis, J. Q., K. Choromanski, J. Varley, H. Lee, J.-J. Slotine, V. Likhosterov, A. Weller, A. Makadia, and V. Sindhvani. (2020). “Time dependence in non-autonomous neural odes”. *arXiv preprint arXiv:2005.01906*.
- DeGroot, M. H. (1962). “Uncertainty, information, and sequential experiments”. *The Annals of Mathematical Statistics*. 33(2): 404–419.
- Dehghani, M., J. Djolonga, B. Mustafa, P. Padlewski, J. Heek, J. Gilmer, A. P. Steiner, M. Caron, R. Geirhos, I. Alabdulmohsin, *et al.* (2023). “Scaling vision transformers to 22 billion parameters”. In: *International Conference on Machine Learning*. PMLR. 7480–7512.
- Deisenroth, M. P., A. A. Faisal, and C. S. Ong. (2020). *Mathematics for machine learning*. Cambridge University Press.
- Devlin, J., M.-W. Chang, K. Lee, and K. Toutanova. (2019). “BERT: Pre-training of deep bidirectional transformers for language understanding”. In: *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*. 4171–4186.
- Drusvyatskiy, D. and C. Paquette. (2019). “Efficiency of minimizing compositions of convex functions and smooth maps”. *Mathematical Programming*. 178: 503–558.
- Duchi, J. C., M. I. Jordan, M. J. Wainwright, and A. Wibisono. (2015). “Optimal rates for zero-order convex optimization: The power of two function evaluations”. *IEEE Transactions on Information Theory*. 61(5): 2788–2806.
- Duchi, J. C., S. Shalev-Shwartz, Y. Singer, and T. Chandra. (2008). “Efficient projections onto the ℓ_1 -ball for learning in high dimensions”. In: *Proc. of ICML*.
- Dufter, P., M. Schmitt, and H. Schütze. (2022). “Position information in transformers: An overview”. *Computational Linguistics*. 48(3): 733–763.
- Eisner, J. (2016). “Inside-outside and forward-backward algorithms are just backprop (tutorial paper)”. In: *Proceedings of the Workshop on Structured Prediction for NLP*. 1–17.

- Elsayed, M. and A. R. Mahmood. (2022). “HesScale: Scalable Computation of Hessian Diagonals”. *arXiv preprint arXiv:2210.11639*.
- Epperly, E. N., J. A. Tropp, and R. J. Webber. (2023). “XTrace: Making the most of every sample in stochastic trace estimation”. *arXiv preprint arXiv:2301.07825*.
- Farinhas, A., W. Aziz, V. Niculae, and A. F. Martins. (2021). “Sparse communication via mixed distributions”. *arXiv preprint arXiv:2108.02658*.
- Flanders, H. (1973). “Differentiation under the integral sign”. *The American Mathematical Monthly*. 80(6): 615–627.
- Fleming, W. H. and R. W. Rishel. (2012). *Deterministic and stochastic optimal control*. Vol. 1. Springer Science & Business Media.
- Fletcher, R. (1970). “A new approach to variable metric algorithms”. *The computer journal*. 13(3): 317–322.
- Foerster, J., G. Farquhar, M. Al-Shedivat, T. Rocktäschel, E. Xing, and S. Whiteson. (2018). “Dice: The infinitely differentiable monte carlo estimator”. In: *International Conference on Machine Learning*. PMLR. 1529–1538.
- Forney, G. D. (1973). “The viterbi algorithm”. *Proceedings of the IEEE*. 61(3): 268–278.
- Franceschi, L., M. Donini, P. Frasconi, and M. Pontil. (2017). “Forward and reverse gradient-based hyperparameter optimization”. In: *International Conference on Machine Learning*. PMLR. 1165–1173.
- Frey, B. J., F. R. Kschischang, H.-A. Loeliger, and N. Wiberg. (1997). “Factor graphs and algorithms”. In: *Proceedings of the Annual Allerton Conference on Communication Control and Computing*. Vol. 35. Citeseer. 666–680.
- Frigyik, B. A., S. Srivastava, and M. R. Gupta. (2008). “An introduction to functional derivatives”. *Dept. Electr. Eng., Univ. Washington, Seattle, WA, Tech. Rep.* 1.
- Frostig, R., M. J. Johnson, D. Maclaurin, A. Paszke, and A. Radul. (2021). “Decomposing reverse-mode automatic differentiation”. *arXiv preprint arXiv:2105.09469*.
- Gautschi, W. (2011). *Numerical analysis*. Springer Science & Business Media.
- Getreuer, P. (2013). “A survey of Gaussian convolution algorithms”. *Image Processing On Line*. 2013: 286–310.

- Geweke, J. (1988). “Antithetic acceleration of Monte Carlo integration in Bayesian inference”. *Journal of Econometrics*. 38(1-2): 73–89.
- Gholaminejad, A., K. Keutzer, and G. Biro. (2019). “ANODE: Unconditionally Accurate Memory-Efficient Gradients for Neural ODEs”. In: *International Joint Conferences on Artificial Intelligence*.
- Gini, C. (1912). “Variabilità e mutabilità”. *Reprinted in Memorie di metodologica statistica (Ed. Pizetti E, Salvemini, T)*. Rome: Libreria Eredi Virgilio Veschi.
- Girard, A. (1989). “A fast ‘Monte-Carlo cross-validation’ procedure for large least squares problems with noisy data”. *Numerische Mathematik*. 56: 1–23.
- Goldfarb, D. (1970). “A family of variable-metric methods derived by variational means”. *Mathematics of computation*. 24(109): 23–26.
- Gomez, A. N., M. Ren, R. Urtasun, and R. B. Grosse. (2017). “The reversible residual network: Backpropagation without storing activations”. *Advances in neural information processing systems*. 30.
- Graves, A., G. Wayne, and I. Danihelka. (2014). “Neural turing machines”. *arXiv preprint arXiv:1410.5401*.
- Greig, D. M., B. T. Porteous, and A. H. Seheult. (1989). “Exact maximum a posteriori estimation for binary images”. *Journal of the Royal Statistical Society Series B: Statistical Methodology*. 51(2): 271–279.
- Griewank, A. (1992). “Achieving logarithmic growth of temporal and spatial complexity in reverse automatic differentiation”. *Optimization Methods and Software*. 1(1): 35–54. DOI: [10.1080/10556789208805505](https://doi.org/10.1080/10556789208805505).
- Griewank, A. (2003). “A mathematical view of automatic differentiation”. *Acta Numerica*. 12: 321–398.
- Griewank, A. (2012). “Who invented the reverse mode of differentiation”. *Documenta Mathematica, Extra Volume ISMP*. 389400.
- Griewank, A. and A. Walther. (2008). *Evaluating derivatives: principles and techniques of algorithmic differentiation*. SIAM.
- Grimm, J., L. Pottier, and N. Rostaing-Schmidt. (1996). “Optimal time and minimum space-time product for reversing a certain class of programs”. *PhD thesis*. INRIA.

- Grünwald, P. D. and A. P. Dawid. (2004). “Game theory, maximum entropy, minimum discrepancy and robust Bayesian decision theory”. *Annals of Statistics*: 1367–1433.
- Hallman, E., I. C. Ipsen, and A. K. Saibaba. (2023). “Monte Carlo methods for estimating the diagonal of a real symmetric matrix”. *SIAM Journal on Matrix Analysis and Applications*. 44(1): 240–269.
- He, K., X. Zhang, S. Ren, and J. Sun. (2016). “Deep residual learning for image recognition”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- Helfrich, K., D. Willmott, and Q. Ye. (2018). “Orthogonal recurrent neural networks with scaled Cayley transform”. In: *International Conference on Machine Learning*. PMLR. 1969–1978.
- Hestenes, M. R., E. Stiefel, *et al.* (1952). *Methods of conjugate gradients for solving linear systems*. Vol. 49. No. 1. NBS Washington, DC.
- Hewitt, E. (1948). “Rings of real-valued continuous functions. I”. *Transactions of the American Mathematical Society*. 64(1): 45–99.
- Hida, T. and M. Hitsuda. (1976). *Gaussian processes*. Vol. 120. American Mathematical Soc.
- Hiriart-Urruty, J.-B. and C. Lemaréchal. (1993). *Convex analysis and minimization algorithms II*. Vol. 305. Springer science & business media.
- Hutchinson, M. F. (1989). “A stochastic estimator of the trace of the influence matrix for Laplacian smoothing splines”. *Communications in Statistics-Simulation and Computation*. 18(3): 1059–1076.
- Imai, T. (2019). *Where did “differentiable programming” come from?* URL: <https://medium.com/@bonotake/where-did-differentiable-programming-come-from-27b385fb6d6d>.
- Ioffe, S. and C. Szegedy. (2015). “Batch normalization: Accelerating deep network training by reducing internal covariate shift”. In: *International conference on machine learning*. pmlr. 448–456.
- Jaggi, M. (2013). “Revisiting Frank-Wolfe: Projection-free sparse convex optimization”. In: *International conference on machine learning*. PMLR. 427–435.
- Jang, E., S. Gu, and B. Poole. (2016). “Categorical reparameterization with gumbel-softmax”. *arXiv preprint arXiv:1611.01144*.

- Jayaram, B. and M. Baczynski. (2008). *Fuzzy Implications*. Vol. 231. Springer Science & Business Media.
- Kakade, S., S. Shalev-Shwartz, A. Tewari, *et al.* (2009). “On the duality of strong convexity and strong smoothness: Learning applications and matrix regularization”. *Tech report*. 2(1): 35.
- Karpathy, A. (2017). “[Software 2.0](#)”.
- Kelley, C. T. (1995). *Iterative methods for linear and nonlinear equations*. SIAM.
- Kingma, D. P. and J. Ba. (2014). “Adam: A method for stochastic optimization”. *arXiv preprint arXiv:1412.6980*.
- Kingma, D. P. and M. Welling. (2013). “Auto-encoding variational bayes”. *arXiv preprint arXiv:1312.6114*.
- Klir, G. and B. Yuan. (1995). *Fuzzy sets and fuzzy logic*. Vol. 4. Prentice hall New Jersey.
- Kobyzev, I., S. Prince, and M. A. Brubaker. (2019). “Normalizing flows: Introduction and ideas”. *stat*. 1050: 25.
- Kreikemeyer, J. N. and P. Andelfinger. (2023). “Smoothing methods for automatic differentiation across conditional branches”. *IEEE Access*.
- Krieken, E., J. Tomczak, and A. Ten Teije. (2021). “Stochastic: A framework for general stochastic automatic differentiation”. *Advances in Neural Information Processing Systems*. 34: 7574–7587.
- Kunstner, F., P. Hennig, and L. Balles. (2019). “Limitations of the empirical Fisher approximation for natural gradient descent”. *Advances in neural information processing systems*. 32.
- Lafferty, J., A. McCallum, and F. C. Pereira. (2001). “Conditional random fields: Probabilistic models for segmenting and labeling sequence data”.
- Lan, G. (2012). “An optimal method for stochastic composite optimization”. *Mathematical Programming*. 133(1-2): 365–397.
- LeCun, Y. (1988). “A theoretical framework for back-propagation”. In: *Proceedings of the 1988 connectionist models summer school*. Vol. 1. 21–28.
- LeCun, Y. (2018). “[Deep Learning est mort. Vive Differentiable Programming!](#)”

- Levenberg, K. (1944). “A method for the solution of certain non-linear problems in least squares”. *Quarterly of applied mathematics*. 2(2): 164–168.
- Liu, D. C. and J. Nocedal. (1989). “On the limited memory method for large scale optimization”. *Mathematical Programming*. 45: 503–528.
- Liu, H., Z. Li, D. Hall, P. Liang, and T. Ma. (2023). “Sophia: A Scalable Stochastic Second-order Optimizer for Language Model Pre-training”. *arXiv preprint arXiv:2305.14342*.
- Loeliger, H.-A. (2004). “An introduction to factor graphs”. *IEEE Signal Processing Magazine*. 21(1): 28–41.
- Loshchilov, I. and F. Hutter. (2016). “SGDR: Stochastic gradient descent with warm restarts”. In: *International Conference on Learning Representations*.
- Lucet, Y. (1997). “Faster than the fast Legendre transform, the linear-time Legendre transform”. *Numerical Algorithms*. 16: 171–185.
- Maclaurin, D., D. Duvenaud, and R. P. Adams. (2015). “Autograd: Effortless gradients in numpy”. In: *ICML 2015 AutoML workshop*. Vol. 238. No. 5.
- Maddison, C. J., A. Mnih, and Y. W. Teh. (2016). “The concrete distribution: A continuous relaxation of discrete random variables”. *arXiv preprint arXiv:1611.00712*.
- Marquardt, D. W. (1963). “An algorithm for least-squares estimation of nonlinear parameters”. *Journal of the society for Industrial and Applied Mathematics*. 11(2): 431–441.
- Martens, J. (2020). “New insights and perspectives on the natural gradient method”. *Journal of Machine Learning Research*. 21(1): 5776–5851.
- Martens, J. and R. Grosse. (2015). “Optimizing neural networks with Kronecker-factored approximate curvature”. In: *International conference on machine learning*. 2408–2417.
- Martins, A. and R. Astudillo. (2016). “From softmax to sparsemax: A sparse model of attention and multi-label classification”. In: *International conference on machine learning*. PMLR. 1614–1623.
- Martins, J. R., P. Sturdza, and J. J. Alonso. (2003). “The complex-step derivative approximation”. *ACM Transactions on Mathematical Software (TOMS)*. 29(3): 245–262.

- Meent, J.-W. van de, B. Paige, H. Yang, and F. Wood. (2018). “An introduction to probabilistic programming”. *arXiv preprint arXiv:1809.10756*.
- Mensch, A. and M. Blondel. (2018). “Differentiable dynamic programming for structured prediction and attention”. In: *International Conference on Machine Learning*. PMLR. 3462–3471.
- Messerer, F., K. Baumgärtner, and M. Diehl. (2021). “Survey of sequential convex programming and generalized Gauss-Newton methods”. *ESAIM: Proceedings and Surveys*. 71: 64–88.
- Meyer, R. A., C. Musco, C. Musco, and D. P. Woodruff. (2021). “Hutch++: Optimal stochastic trace estimation”. In: *Symposium on Simplicity in Algorithms (SOSA)*. SIAM. 142–155.
- Michelot, C. (1986). “A finite algorithm for finding the projection of a point onto the canonical simplex of \mathbb{R}^n ”. *Journal of Optimization Theory and Applications*. 50(1): 195–200.
- Mohamed, S., M. Rosca, M. Figurnov, and A. Mnih. (2020). “Monte carlo gradient estimation in machine learning”. *The Journal of Machine Learning Research*. 21(1): 5183–5244.
- Mohri, M., F. Pereira, and M. Riley. (2008). “Speech recognition with weighted finite-state transducers”. *Springer Handbook of Speech Processing*: 559–584.
- Morgenstern, J. (1985). “How to compute fast a function and all its derivatives: A variation on the theorem of Baur-Strassen”. *ACM SIGACT News*. 16(4): 60–62.
- Morrey Jr, C. B. (2009). *Multiple integrals in the calculus of variations*. Springer Science & Business Media.
- Murphy, K. P. (2022). *Probabilistic Machine Learning: An introduction*. MIT Press. URL: <http://probml.github.io/book1>.
- Murphy, K. P. (2023). *Probabilistic Machine Learning: Advanced Topics*. MIT Press. URL: <http://probml.github.io/book2>.
- Mutze, U. (2013). “An asynchronous leapfrog method II”. *arXiv preprint arXiv:1311.6602*.
- Nemirovski, A. and D. Yudin. (1983). “Problem complexity and method efficiency in optimization”.
- Nesterov, Y. (2005). “Smooth minimization of non-smooth functions”. *Mathematical programming*. 103: 127–152.

- Nesterov, Y. (2007). “Modified Gauss–Newton scheme with worst case guarantees for global performance”. *Optimisation methods and software*. 22(3): 469–483.
- Nesterov, Y. (2018). *Lectures on convex optimization*. Vol. 137. Springer.
- Nesterov, Y. and V. Spokoiny. (2017). “Random gradient-free minimization of convex functions”. *Foundations of Computational Mathematics*. 17: 527–566.
- Olah, C. (2015). *Neural networks, types, and functional programming*. URL: <https://colah.github.io/posts/2015-09-NN-Types-FP/>.
- Papamakarios, G., E. Nalisnick, D. J. Rezende, S. Mohamed, and B. Lakshminarayanan. (2021). “Normalizing flows for probabilistic modeling and inference”. *The Journal of Machine Learning Research*. 22(1): 2617–2680.
- Parikh, N., S. Boyd, *et al.* (2014). “Proximal algorithms”. *Foundations and trends® in Optimization*. 1(3): 127–239.
- Paszke, A., S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. (2019). “PyTorch: An Imperative Style, High-Performance Deep Learning Library”. In: *Advances in Neural Information Processing Systems 32*. 8024–8035.
- Paulus, M., D. Choi, D. Tarlow, A. Krause, and C. J. Maddison. (2020). “Gradient estimation with stochastic softmax tricks”. *Advances in Neural Information Processing Systems*. 33: 5691–5704.
- Petersen, F., C. Borgelt, H. Kuehne, and O. Deussen. (2021). “Learning with algorithmic supervision via continuous relaxations”. *Advances in Neural Information Processing Systems*. 34: 16520–16531.
- Peyré, G. (2020). “Mathematical foundations of data sciences”. *Rn*. 1: 2.
- Peyré, G. and M. Cuturi. (2019). “Computational optimal transport: With applications to data science”. *Foundations and Trends® in Machine Learning*. 11(5-6): 355–607.
- Plotkin, G. (2018). *Some Principles of Differential Programming Languages*. URL: <https://popl18.sigplan.org/details/POPL-2018-papers/76/Some-Principles-of-Differential-Programming-Languages>.

- Pollock, S. and L. G. Rebholz. (2021). “Anderson acceleration for contractive and noncontractive operators”. *IMA Journal of Numerical Analysis*. 41(4): 2841–2872.
- Polyak, B. (1963). “Gradient methods for the minimisation of functionals”. *USSR Computational Mathematics and Mathematical Physics*. 3(4): 864–878.
- Polyak, B. T. (1964). “Some methods of speeding up the convergence of iteration methods”. *Ussr computational mathematics and mathematical physics*. 4(5): 1–17.
- Pontryagin, L. S. (1985). “The mathematical theory of optimal processes and differential games”. *Trudy Mat. Inst. Steklov*. 169: 119–158.
- Press, O. and L. Wolf. (2016). “Using the output embedding to improve language models”. *arXiv preprint arXiv:1608.05859*.
- Rabiner, L. R. (1989). “A tutorial on hidden Markov models and selected applications in speech recognition”. *Proceedings of the IEEE*. 77(2): 257–286.
- Rademacher, H. (1919). “Über partielle und totale differenzierbarkeit von Funktionen mehrerer Variabeln und über die Transformation der Doppelintegrale”. *Mathematische Annalen*. 79(4): 340–359.
- Radul, A., A. Paszke, R. Frostig, M. Johnson, and D. Maclaurin. (2022). “You only linearize once: Tangents transpose to gradients”. *arXiv preprint arXiv:2204.10923*.
- Rahimi, A. and B. Recht. (2007). “Random features for large-scale kernel machines”. *Advances in neural information processing systems*. 20.
- Recht, B. (2016). “Mates of Costate”.
- Recht, B. and R. Frostig. (2017). “Nesterov’s Punctuated Equilibrium”.
- Rezende, D. J., S. Mohamed, and D. Wierstra. (2014). “Stochastic back-propagation and approximate inference in deep generative models”. In: *International conference on machine learning*. PMLR. 1278–1286.
- Rockafellar, R. T. and R. J.-B. Wets. (2009). *Variational analysis*. Vol. 317. Springer Science & Business Media.
- Rodriguez, O. H. and J. M. Lopez Fernandez. (2010). “A semiotic reflection on the didactics of the chain rule”. *The Mathematics Enthusiast*. 7(2): 321–332.

- Roulet, V. and Z. Harchaoui. (2022). “Differentiable programming à la Moreau”. In: *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE. 3498–3502.
- Saad, Y. and M. H. Schultz. (1986). “GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems”. *SIAM Journal on scientific and statistical computing*. 7(3): 856–869.
- Salimans, T., J. Ho, X. Chen, S. Sidor, and I. Sutskever. (2017). “Evolution strategies as a scalable alternative to reinforcement learning”. *arXiv preprint arXiv:1703.03864*.
- Sander, M. E., P. Ablin, M. Blondel, and G. Peyré. (2021a). “Momentum residual neural networks”. In: *International Conference on Machine Learning*. PMLR. 9276–9287.
- Sander, M. E., P. Ablin, M. Blondel, and G. Peyré. (2021b). “Momentum residual neural networks”. In: *International Conference on Machine Learning*. PMLR. 9276–9287.
- Satterthwaite, F. (1942). “Generalized poisson distribution”. *The Annals of Mathematical Statistics*. 13(4): 410–417.
- Schlag, I., K. Irie, and J. Schmidhuber. (2021). “Linear transformers are secretly fast weight programmers”. In: *International Conference on Machine Learning*. PMLR. 9355–9366.
- Schölkopf, B. and A. J. Smola. (2002). *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press.
- Schulman, J., N. Heess, T. Weber, and P. Abbeel. (2015). “Gradient estimation using stochastic computation graphs”. *Advances in neural information processing systems*. 28.
- Schwartz, J. (1954). “The formula for change in variables in a multiple integral”. *The American Mathematical Monthly*. 61(2): 81–85.
- Schwarz, H. (1873). “Communication”. *Archives des Sciences Physiques et Naturelles*. 48: 38–44.
- Sengupta, S., M. J. Harris, M. Garland, and J. D. Owens. (2010). “Efficient Parallel Scan Algorithms for Manycore GPUs.”
- Shanno, D. F. (1970). “Conditioning of quasi-Newton methods for function minimization”. *Mathematics of computation*. 24(111): 647–656.

- Shannon, C. E. (1948). “A mathematical theory of communication”. *The Bell system technical journal*. 27(3): 379–423.
- Shawe-Taylor, J. and N. Cristianini. (2004). *Kernel methods for pattern analysis*. Cambridge university press.
- Squire, W. and G. Trapp. (1998). “Using complex variables to estimate derivatives of real functions”. *SIAM review*. 40(1): 110–112.
- Stoer, J., R. Bulirsch, R. Bartels, W. Gautschi, and C. Witzgall. (1980). *Introduction to numerical analysis*. Vol. 1993. Springer.
- Stumm, P. and A. Walther. (2010). “New algorithms for optimal online checkpointing”. *SIAM Journal on Scientific Computing*. 32(2): 836–854.
- Su, J., M. Ahmed, Y. Lu, S. Pan, W. Bo, and Y. Liu. (2024). “Ro-former: Enhanced transformer with rotary position embedding”. *Neurocomputing*. 568: 127063.
- Sutherland, D. J. and J. Schneider. (2015). “On the error of random Fourier features”. *arXiv preprint arXiv:1506.02785*.
- Sutskever, I., J. Martens, G. Dahl, and G. Hinton. (2013). “On the importance of initialization and momentum in deep learning”. In: *International conference on machine learning*. PMLR. 1139–1147.
- Sutton, C., A. McCallum, et al. (2012). “An introduction to conditional random fields”. *Foundations and Trends® in Machine Learning*. 4(4): 267–373.
- Sutton, R. S., D. McAllester, S. Singh, and Y. Mansour. (1999). “Policy gradient methods for reinforcement learning with function approximation”. *Advances in neural information processing systems*. 12.
- Taylor, M. (2002). “Differential forms and the change of variable formula for multiple integrals”. *Journal of mathematical analysis and applications*. 268(1): 378–383.
- Tibshirani, R. (1996). “Regression shrinkage and selection via the lasso”. *Journal of the Royal Statistical Society: Series B (Methodological)*. 58(1): 267–288.
- Tignol, J.-P. (2015). *Galois’ theory of algebraic equations*. World Scientific Publishing Company.
- Tsallis, C. (1988). “Possible generalization of Boltzmann-Gibbs statistics”. *Journal of statistical physics*. 52: 479–487.

- van Krieken, E. (2024). “Optimisation in Neurosymbolic Learning Systems”. *PhD thesis*. Vrije Universiteit Amsterdam.
- Vaswani, A., N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. (2017). “Attention is all you need”. *Advances in neural information processing systems*. 30.
- Vaswani, S., A. Mishkin, I. Laradji, M. Schmidt, G. Gidel, and S. Lacoste-Julien. (2019). “Painless stochastic gradient: Interpolation, line-search, and convergence rates”. *Advances in neural information processing systems*. 32.
- Verdu, S. and H. V. Poor. (1987). “Abstract dynamic programming models under commutativity conditions”. *SIAM Journal on Control and Optimization*. 25(4): 990–1006.
- Vicol, P., L. Metz, and J. Sohl-Dickstein. (2021). “Unbiased gradient estimation in unrolled computation graphs with persistent evolution strategies”. In: *International Conference on Machine Learning*. PMLR. 10553–10563.
- Virtanen, P., R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, Í. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors. (2020). “SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python”. *Nature Methods*. 17: 261–272.
- Viterbi, A. (1967). “Error bounds for convolutional codes and an asymptotically optimum decoding algorithm”. *IEEE transactions on Information Theory*. 13(2): 260–269.
- Vorst, H. A. v. d. and H. A. van der Vorst. (1992). “Bi-CGSTAB: A Fast and Smoothly Converging Variant of Bi-CG for the Solution of Nonsymmetric Linear Systems”. *SIAM Journal on Scientific and Statistical Computing*. 13(2): 631–644. URL: <http://dx.doi.org/10.1137/0913035>.
- Wainwright, M. J. and M. I. Jordan. (2008). “Graphical models, exponential families, and variational inference”. *Foundations and Trends® in Machine Learning*. 1(1–2): 1–305.

- Wang, Q., P. Moin, and G. Iaccarino. (2009). “Minimal repetition dynamic checkpointing algorithm for unsteady adjoint calculation”. *SIAM Journal on Scientific Computing*. 31(4): 2549–2567.
- Wei, C., S. Kakade, and T. Ma. (2020). “The implicit and explicit regularization effects of dropout”. In: *International conference on machine learning*. PMLR. 10181–10192.
- Werbos, P. J. (1990). “Backpropagation through time: what it does and how to do it”. *Proceedings of the IEEE*. 78(10): 1550–1560.
- Werbos, P. J. (1994). *The roots of backpropagation: from ordered derivatives to neural networks and political forecasting*. Vol. 1. John Wiley & Sons.
- Wright, S. and J. Nocedal. (1999). “Numerical optimization”. *Springer Science*. 35(67-68): 7.
- Xu, P., F. Roosta, and M. W. Mahoney. (2020). “Second-order optimization for non-convex machine learning: An empirical study”. In: *Proceedings of the 2020 SIAM International Conference on Data Mining*. SIAM. 199–207.
- Yuan, M. and Y. Lin. (2006). “Model selection and estimation in regression with grouped variables”. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 68(1): 49–67.
- Zhang, A., Z. C. Lipton, M. Li, and A. J. Smola. (2021). “Dive into deep learning”. *arXiv preprint arXiv:2106.11342*.
- Zhou, X. (2018). “On the fenchel duality between strong convexity and lipschitz continuous gradient”. *arXiv preprint arXiv:1803.06573*.
- Zhuang, J., N. C. Dvornik, S. Tatikonda, and J. S. Duncan. (2021). “Mali: A memory efficient and reverse accurate integrator for neural odes”. *arXiv preprint arXiv:2102.04668*.
- Ziegler, D. M., N. Stiennon, J. Wu, T. B. Brown, A. Radford, D. Amodei, P. Christiano, and G. Irving. (2019). “Fine-tuning language models from human preferences”. *arXiv preprint arXiv:1909.08593*.