



Statistical Inference and Learning

Author: Wenxiao Yang

Institute: Haas School of Business, University of California Berkeley

Date: 2023

All models are wrong, but some are useful.

Contents

Chapter 1 Information-Theoretic Functional	1
1.1 Definitions	1
1.1.1 Entropy	1
1.1.2 Kullback-Leibler Divergence	1
1.1.3 Cross-Entropy	2
1.1.4 Mutual Information	2
Chapter 2 Statistics Basics	3
2.1 Random Sampling	3
2.1.1 Sample Mean and Sample Variance	3
2.1.2 Distributional Properties	4
2.1.3 Order Statistics	4
2.2 Basic Statistics	5
2.3 Point Estimation	6
2.3.1 Method of Moments (MM)	6
2.3.2 Maximum Likelihood (ML)	8
2.3.3 Comparing Estimators: Mean Squared Error	9
2.3.4 Sufficient Statistic	10
2.3.5 Minimal Sufficient Statistic	12
2.3.6 Complete Statistic	13
2.3.7 Cramér-Rao Lower Bound	14
2.4 Hypothesis Testing	16
2.4.1 Formulation of Testing Problem	16
2.4.2 Errors, Power Function, and Agenda	17
2.4.3 Choice of Critical Value	18
2.4.4 Choice of Test Statistic: Uniformly Most Powerful (UMP) Level α Test	19
2.4.5 Generalized Neyman-Pearson Lemma	21
2.5 Trinity of Classical Tests	21
2.5.1 Test Statistics	22

2.5.2	Approximation to T_{LR}	23
2.6	Interval Estimation	23
Chapter 3 Decision Rule Based Statistical Inference		25
3.1	Decision Rule	25
3.2	Maximum-Likelihood Principle (state is norandom)	25
3.3	Bayesian Decision Rule (state is random)	26
3.3.1	Rules	26
3.3.2	Optimization Problem in Bayes Form	27
3.3.3	Maximum A Posteriori (MAP) Decision Rule (Binary example)	29
3.3.4	Minimum Mean Squared Error (MMSE) Rule (\mathbb{R}^n example)	29
3.4	Comparison	30
Chapter 4 Bootstrap		31
4.1	Traditional Monte-Carlo Approach	31
4.2	Bootstrap (When data is not enough)	32
4.3	Residual Bootstrap (for problem with not i.i.d. data)	32
4.3.1	Example: Linear	33
4.3.2	Example: Nonlinear Markov Process	33
4.4	Posterior Simulation / Bayesian (Weighted) Bootstrap	34
4.4.1	Dirichlet Distribution Prior	34
4.4.2	Haldane Prior	35
4.4.3	Linear Model Case	35
4.4.4	Bernoulli Case	36
Chapter 5 Nonparameteric Prediction Probelm		37
5.1	K -normal Means Probelm	38
5.1.1	Assumptions	38
5.1.2	Maximum Likelihood Estimator	39
5.1.3	Risk of MLE	39
5.1.4	James-Stein Type Estimator	40
Chapter 6 Linear Predictors / Regression		43
6.1	Best Linear Predictor	43

6.2	Convergence of OLS	44
6.2.1	Approximation	44
6.2.2	Testing and Confidence Interval	46
6.3	Long, Short, Auxiliary Regression	46
6.4	Residual Regression	48
6.5	Card-Krueger Model	49
6.5.1	Proxy Variable Regression	50
Chapter 7	Machine Learning in Inference	52
7.1	Empirical Risk Minimization (ERM)	52
7.1.1	Example: Linear MMSE (LMMSE) estimator	52
7.1.2	Penalized ERM	53
7.2	Stochastic Approximation	54
7.3	Stochastic Gradient Descent (SGD)	56
7.4	SGD Application to Empirical Risk Minimization (ERM)	57
7.4.1	Different Gradient Descent for ERM	58
7.4.2	Constraints on Learning Problem	58
Chapter 8	Stochastic Integration Methods	60
8.1	Deterministic Methods (Better in Low Dimension)	60
8.1.1	Riemann Integration	60
8.1.2	Trapezoidal Rule	60
8.1.3	Multidimensional Integration	61
8.2	Stochastic Methods (Better in High Dimension)	61
8.2.1	Classical Monte Carlo Integration	61
8.2.2	Importance Sampling	62
Chapter 9	Particle Filtering	64
9.1	Kalman Filtering (Linear Dynamic System)	64
9.2	Particle Filtering (Nonlinear Dynamic System)	64
9.2.1	Bayesian Recursive Filtering	65
9.2.2	Particle Filter (bootstrap filter)	65
Chapter 10	EM Algorithm	67

10.1 General Structure of the EM Algorithm	67
10.2 Example 1: Variance Estimation	69
10.2.1 Maximum-Likelihood (ML) Estimation	69
10.2.2 EM Algorithm	69
10.3 Example 2: Estimation of Gaussian Mixtures	70
10.3.1 Unknown Means: ML estimation is hard	70
10.3.2 Unknown Means: EM Algorithm	71
10.3.3 Unknown Mixture Probabilities, Means and Variances	72
10.4 Convergence of EM Algorithm	72
10.5 EM As an Alternating Maximization Algorithm	73
Chapter 11 Hidden Markov model (HMM)	75
11.1 Viterbi Algorithm: (MAP) estimate $X_{1:t}$ given $Y_{1:t}$	75
11.1.1 MAP estimation problem	75
11.1.2 Viterbi Algorithm	76
11.2 Bayesian Estimation of a Sequence: Need (MMSE) estimate $X_{1:t}$ given $Y_{1:t}$	77
11.3 Forward-Backward Algorithm: (MMSE) estimate $X_{1:t+1}$ given $Y_{1:t}$	77
11.3.1 $\gamma_t(x) \triangleq P\left\{X_t = x \mid \vec{Y} = \vec{y}\right\}$	77
11.3.2 $\xi_t(x, x') \triangleq P\left\{X_t = x, X_{t+1} = x' \mid \vec{Y} = \vec{y}\right\}$	79
11.3.3 Scaling Factors	79
Chapter 12 Graphic Models	80
12.1 Graph Theory	80
12.2 Bayesian Networks	81
12.3 Markov Networks	81
12.3.1 General Form	81
12.3.2 Hammersley-Clifford theorem	82
12.3.3 Form of Gibbs distribution (Boltzmann distribution)	82
12.4 Conversion of directed graph to undirected graph	83
12.5 Inference and Learning	83
12.5.1 Inference on Trees	83
Chapter 13 Variational Inference, Mean-Field Techniques	86
13.1 Naive Mean-Field Methods	86

13.1.1 Graphical Models	87
13.1.2 Ising Model	87
13.2 Exponential Families of Probability Distributions	89
13.3 ML Estimation	91
13.4 Maximum Entropy	91
13.5	92
13.6 Connection between Exponential Families and Graphic Models	93
13.6.1 Marginal polytope	93
13.6.2 Locally Consistent Marginal Distributions	93
13.6.3 Entropy on Tree Graphs	95
13.6.4 Naive Mean-Field Methods In Graph	96
13.6.5 Structural Mean Field Optimization	96
13.6.6 Bethe Entropy Approximation	96
Chapter 14 ℓ_1 Penalized Least Squares Minimization	98
14.1 Problem Statement	98
14.2 Special Cases	99
14.2.1 Definition: Soft Threshold	99
14.2.2 Identity A	99
14.2.3 Orthonormal A	99
14.2.4 Quadratic Optimization ($\lambda = 0$)	99
14.3 General Solution: Lasso	100
14.4 General Solution: Iterative Soft Thresholding Algorithm (ISTA)	100
14.4.1 Proximal Minimization Algorithm	100
14.4.2 Apply to ℓ_1 -penalized least-squares	101
14.5 Convergence Rate	101
14.6 Fast Iterative Soft Thresholding Algorithm (FISTA)	102
14.7 Alternating Direction Method of Multipliers (ADMM)	102
Chapter 15 Compressive Sensing	104
15.1 Definitions related to Sparsity	104
15.2 Measurement Matrix	106
15.2.1 Matrix Preliminaries	106

15.2.2 Recovery of k-Sparse Signals	107
15.2.3 Restricted Isometry Property	108
15.3 Robust Signal Recovery from Noiseless Observations	108
15.4 Robust Signal Recovery from Noisy Observations	110
15.4.1 Bounded Noise	110

Chapter 1 Information-Theoretic Functional

1.1 Definitions

1.1.1 Entropy

Definition 1.1 (Entropy)

Entropy of pmf $\{p(x), x \in X\}$

$$H(p) = - \sum_{x \in X} p(x) \ln p(x)$$

(concave in p)



1.1.2 Kullback-Leibler Divergence

Definition 1.2 (KL divergence, Relative entropy)

The **Kullback-Leibler divergence** (or relative entropy) of two pmf's $p(x), x \in X$ and $q(x), x \in X$ is defined as

$$D(p\|q) = - \sum_{x \in X} p(x) \ln \frac{q(x)}{p(x)}$$



Proposition 1.1 (Positivity)

$D(p\|q) = - \sum_{x \in X} p(x) \ln \frac{q(x)}{p(x)} \geq 0$ with equality iff $p = q$. (convex in (p, q))



Proof 1.1 (Using Jensen's Inequality)

Let $g(X) = \frac{q(X)}{p(X)}$, then $D(p\|q)$ can be written as $D(p\|q) = - \sum_{x \in X} p(x) \ln g(x)$.

Because \log is concave, by the Jensen's inequality, we have

$$\sum_{x \in X} p(x) \ln g(x) \leq \ln \sum_{x \in X} p(x)g(x) = \ln \sum_{x \in X} q(X) = \ln 1 = 0$$

where the inequality achieves equality if and only if $g(x)$ is constant for all $x \in X$. (i.e., $g(x) = \int_{-\infty}^{\infty} p(x)g(x)dx = 1$)

Proof 1.2 (Alternative proof by $\ln x \leq x - 1$)

Because $\ln x \leq x - 1$ for all $x > 0$,

$$D(p\|q) = - \sum_{x \in X} p(x) \ln \frac{q(x)}{p(x)} \geq - \sum_{x \in X} p(x) \left(\frac{q(x)}{p(x)} - 1 \right) = - \sum_{x \in X} (q(x) - p(x)) = 0$$

where the inequality achieves equality if and only if $\frac{q(x)}{p(x)} = 1$.

1.1.3 Cross-Entropy

Definition 1.3 (Cross-entropy)

The **cross-entropy** of a pmf $p(x), x \in X$ relative to another pmf $q(x), x \in X$

$$\begin{aligned} H(p, q) &= - \sum_{x \in X} p(x) \ln q(x) \\ &= H(p) + D(p\|q) \end{aligned}$$

$H(p, q) \geq H(p)$, the lower bound is achieved by $q = p$.



1.1.4 Mutual Information

Definition 1.4 (Mutual Information)

Let (x, y) be a pair of random variables with values over the space $X \times Y$. If their joint distribution is $p_{X,Y}(x, y)$ and the marginal distributions are $p_X(x)$ and $p_Y(y)$, the **mutual information** is defined as

$$\begin{aligned} I(p_{X,Y}) &= \sum_{x \in X} \sum_{y \in Y} p_{X,Y}(x, y) \log \left(\frac{p_{X,Y}(x, y)}{p_X(x)p_Y(y)} \right) \\ &= H(p_X) + H(p_Y) - H(p_X, p_Y) \end{aligned}$$



Chapter 2 Statistics Basics

2.1 Random Sampling

Definition 2.1 (Random Sample)

A **random sample** is a collection X_1, \dots, X_n of random variables that are (mutually) independent and identical marginal distributions.

X_1, \dots, X_n are called "independent and identically distributed". The notation is $X_i \sim i.i.d.$



Definition 2.2 (Statistic)

A **statistic** (singular) or sample statistic is any quantity computed from values in a sample which is considered for a statistical purpose.

If X_1, \dots, X_n is a random sample and $T : \mathbb{R}^n \rightarrow \mathbb{R}^k$ (for some $k \in \mathbb{N}$), then $T(X_1, \dots, X_n)$ is called a **statistic**.



2.1.1 Sample Mean and Sample Variance

Definition 2.3 (Sample Mean and Sample Variance)

1. The **sample mean** is $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$;
2. The **sample variance** is $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} (\sum_{i=1}^n X_i^2 - n\bar{X}^2)$



Note We use " $X_i \sim i.i.d(\mu, \sigma^2)$ " to denote a random sample from a distribution with mean μ and variance σ^2 .

Theorem 2.1 ($\mathbb{E}(\bar{X}), \text{Var}(\bar{X}), \mathbb{E}(S^2)$)

Suppose X_1, \dots, X_n is a random sample from a distribution with mean μ and variance σ^2 (denoted by $X_i \sim i.i.d(\mu, \sigma^2)$). Then,

- (a). $\mathbb{E}(\bar{X}) = \mu$;
- (b). $\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$;
- (c). $\mathbb{E}(S^2) = \sigma^2$.



2.1.2 Distributional Properties

Theorem 2.2

If $X_i \sim \text{i.i.d. } N(\mu, \sigma^2)$, then

- (a). $\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$
- (b). $\frac{n-1}{\sigma^2} S^2 \sim \chi_{n-1}^2$
- (c). $\bar{X} \perp S^2$



Theorem 2.3 ("Asymptotics")

If $X_i \sim \text{i.i.d. } (\mu, \sigma^2)$ and if n is "large", then

- (a). $\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$ (converges in distribution) by CLT 6.2;
- (b). $S^2 = \sigma^2$ by LLN;



2.1.3 Order Statistics

Definition 2.4 (Order Statistics)

If X_1, \dots, X_n is a random sample, then the **characteristics** are the sample values placed in ascending order. Notation:

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$$



Proposition 2.1 (Distribution of $X_n = \max_{i=1,\dots,n} X_i$)

If X_1, \dots, X_n is a random sample from a distribution with cdf F (denoted by " $X_i \sim \text{i.i.d. } F$ "), then

$$F_{X_{(n)}}(x) = P(X_{(n)} \leq x) = F^n(x)$$



Proposition 2.2 (cdf and pdf)

More generally,

$$\begin{aligned} F_{X_{(r)}}(x) &= \sum_{j=r}^n \binom{n}{j} [F_X(x)]^j [1 - F_X(x)]^{n-j} \\ f_{X_{(r)}}(x) &= \frac{n!}{(r-1)!(n-r)!} f_X(x) [F_X(x)]^{r-1} [1 - F_X(x)]^{n-r} \end{aligned}$$



Example 2.1

1. **Order statistics sampled from a uniform distribution on unit interval ($\text{Unif}[0, 1]$):** Consider a random sample U_1, \dots, U_n from the standard uniform distribution. Then,

$$f_{X_{(k)}}(x) = \frac{n!}{(k-1)!(n-k)!} u^{k-1} (1-u)^{n-k}$$

The k^{th} order statistic of the uniform distribution is a beta-distributed random variable.

$$U_{(k)} \sim \text{Beta}(k, n+1-k)$$

which has mean $\mathbb{E}[U_{(k)}] = \frac{k}{n+1}$.

2. The joint distribution of the order statistics of the uniform distribution on unit interval ($\text{Unif}[0, 1]$):

Similarly, for $i < j$, the joint probability density function of the two order statistics $U_{(i)} < U_{(j)}$ can be shown to be

$$f_{U_{(i)}, U_{(j)}}(u, v) = n! \frac{u^{i-1}}{(i-1)!} \frac{(v-u)^{j-i-1}}{(j-i-1)!} \frac{(1-v)^{n-j}}{(n-j)!}$$

The joint density of the n order statistics turns out to be constant:

$$f_{U_{(1)}, U_{(2)}, \dots, U_{(n)}}(u_1, u_2, \dots, u_n) = n!$$

For $n \geq k > j \geq 1$, $U_{(k)} - U_{(j)}$ also has a beta distribution:

$$U_{(k)} - U_{(j)} \sim \text{Beta}(k-j, n-(k-j)+1)$$

which has mean $\mathbb{E}[U_{(k)} - U_{(j)}] = \frac{k-j}{n+1}$

2.2 Basic Statistics

In statistics, we define **data** be a vector $x = (x_1, \dots, x_n)'$ of numbers.

Assumption [Fundamental Assumption] x is the realization of a random vector $X = (X_1, \dots, X_n)'$.

Objective: Using x to give (data-based) answers to questions about the distribution of X .

Probability vs. Statistics:

- Probability: Distribution known, outcome unknown;
- Statistics: Distribution unknown, outcome known.

Setting: X_1, \dots, X_n is a random sample from a discrete/continuous distribution with pmf/pdf $f(\cdot | \theta)$, where $\theta \in \Theta$ is unknown.

Types of Statistical Inference:

- Point estimation \Rightarrow "What is θ ?";
- Hypothesis testing \Rightarrow "Is $\theta = \theta_0$?";
- Interval estimation \Rightarrow "Which values of θ are 'plausible'?".

Example 2.2 Examples of Statistical Models

- (1). $x_i \sim \text{i.i.d. Bernoulli}(p)$, where p is unknown.
- (2). $x_i \sim \text{i.i.d. } U(0, \theta)$, where $\theta > 0$ is unknown.
- (3). $x_i \sim \text{i.i.d. } N(\mu, \sigma^2)$, where $\mu \in \mathbb{R}$ and $\sigma^2 > 0$ are unknown.

2.3 Point Estimation

Suppose X_1, \dots, X_n is a random sample from a discrete/continuous distribution with pmf/pdf $f(\cdot | \theta)$, where $\theta \in \Theta$ is unknown.

Definition 2.5 (Point Estimator)

A **point estimator** (of θ) is a function of (X_1, \dots, X_n) .

Notation: $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$.



Agenda

(1). Constructing point estimators

- o Method of moments;
- o Maximum likelihood.

(2). Comparing estimators

- o Pairwise comparisons;
- o Finding 'optimal' estimators.

2.3.1 Method of Moments (MM)

Definition 2.6 (Method of Moments in \mathbb{R}^1)

Suppose $\Theta \subseteq \mathbb{R}^1$. A **method of moments** estimator $\hat{\theta}_{MM}$ solves

$$\mu(\hat{\theta}_{MM}) = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

where $\mu : \Theta \rightarrow \mathbb{R}$ is given by

$$\mu(\theta) = \begin{cases} \sum_{x \in \mathbb{R}} xf(x | \theta), & \text{if } X_i \text{ are discrete} \\ \int_{-\infty}^{\infty} xf(x | \theta) dx, & \text{if } X_i \text{ are continuous} \end{cases}$$



Remark Existence of $\mu(\cdot)$ is assumed; Existence (and uniqueness) of $\hat{\theta}_{MM}$ is assumed.

Example 2.3

1. Suppose $X_i \sim$ i.i.d. $Ber(p)$ where $p \in [0, 1]$ is unknown. The moment function is

$$\mu(p) = p$$

Then, the estimator is

$$\hat{p}_{MM} = \mu(\hat{p}_{MM}) = \bar{X}$$

Remark $\hat{p}_{MM} = \bar{X}$ is the 'best' estimator of p .

2. Suppose $X_i \sim \text{i.i.d.}U(0, \theta)$ where $\theta > 0$ is unknown.

Remark Non-regular statistical model: parameter dependent support, where $\text{supp } X = [0, \theta]$.

The moment function is

$$\mu(\theta) = \frac{\theta}{2}$$

Then, the estimator is

$$\hat{\theta}_{MM} = 2\mu(\hat{\theta}_{MM}) = 2\bar{X}$$

Remark $\hat{\theta}_{MM}$ is not a very good estimator of θ . Concern $X_i > \hat{\theta}_{MM}$ could happen. So, $\max\{\hat{\theta}_{MM}, X_{(n)}\}$ can be better.

Definition 2.7 (Method of Moments in \mathbb{R}^k)

Suppose $\Theta \subseteq \mathbb{R}^k$. A **method of moments** estimator $\hat{\theta}_{MM}$ solves

$$\mu'_j(\hat{\theta}_{MM}) = \frac{1}{n} \sum_{i=1}^n X_i^j, \quad (j = 1, \dots, k)$$

where $\mu'_j : \Theta \rightarrow \mathbb{R}$ is given by

$$\mu'_j(\theta) = \begin{cases} \sum_{x \in \mathbb{R}} x^j f(x | \theta), & \text{if } X_i \text{ are discrete} \\ \int_{-\infty}^{\infty} x^j f(x | \theta) dx, & \text{if } X_i \text{ are continuous} \end{cases}$$



Example 2.4

Suppose $X_i \sim \text{i.i.d.}N(\mu, \sigma^2)$ where $\mu \in \mathbb{R}$ and $\sigma^2 > 0$ are unknown. The moment function is

$$\mu'_1(\mu, \sigma^2) = \mu$$

$$\mu'_2(\mu, \sigma^2) = \mu^2 + \sigma^2$$

Then, the estimator is

$$\begin{aligned} \mu'_1(\hat{\mu}_{MM}, \hat{\sigma}_{MM}^2) &= \hat{\mu}_{MM} = \frac{1}{n} \sum_{i=1}^n X_i \\ \mu'_2(\hat{\mu}_{MM}, \hat{\sigma}_{MM}^2) &= \hat{\mu}_{MM} + \hat{\sigma}_{MM}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 \\ \Rightarrow \hat{\mu}_{MM} &= \bar{X} \\ \hat{\sigma}_{MM}^2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \end{aligned}$$

Remark \bar{X} is the 'best' estimator of μ ; An alternative better estimator of σ^2 is $\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$.

2.3.2 Maximum Likelihood (ML)

Definition 2.8 (Maximum Likelihood)

A **maximum likelihood estimator** $\hat{\theta}_{ML}$ solves

$$L(\hat{\theta}_{ML} \mid X_1, \dots, X_n) = \max_{\theta \in \Theta} L(\theta \mid X_1, \dots, X_n)$$

where $L(\cdot \mid X_1, \dots, X_n) : \Theta \rightarrow \mathbb{R}_+$ is given by

$$L(\theta \mid X_1, \dots, X_n) = \prod_{i=1}^n f_{X_i}(X_i \mid \theta), \quad \theta \in \Theta$$



Remark $L(\cdot \mid X_1, \dots, X_n)$ is called the likelihood function.

Definition 2.9 (Log-Likelihood)

The **log-likelihood** function is

$$l(\theta \mid X_1, \dots, X_n) = \log L(\theta \mid X_1, \dots, X_n) = \sum_{i=1}^n \log f_{X_i}(X_i \mid \theta), \quad \theta \in \Theta$$



Example 2.5

- Suppose $X_i \sim \text{i.i.d. Ber}(p)$ where $p \in [0, 1]$ is unknown. The marginal pmf is

$$f(x \mid p) = \begin{cases} p, & x = 1 \\ 1 - p, & x = 0 \\ 0, & \text{otherwise} \end{cases} = p^x (1 - p)^{1-x} \mathbf{1}_{\{x \in \{0,1\}\}}$$

Then, the likelihood function is

$$\begin{aligned} L(p \mid X_1, \dots, X_n) &= \prod_{i=1}^n \left\{ p^{X_i} (1 - p)^{1-X_i} \underbrace{\mathbf{1}_{\{X_i \in \{0,1\}\}}}_{=1} \right\} \\ &= p^{\sum_{i=1}^n X_i} (1 - p)^{n - \sum_{i=1}^n X_i}, \quad p \in [0, 1] \end{aligned}$$

and the log-likelihood function is

$$l(p \mid X_1, \dots, X_n) = \left(\sum_{i=1}^n X_i \right) \log p + \left(n - \sum_{i=1}^n X_i \right) \log(1 - p), \quad p \in (0, 1)$$

Maximization:

- (a). Suppose $0 < \sum_{i=1}^n X_i < n$, we can give the first-order condition:

$$\begin{aligned} \frac{\partial l(p \mid X_1, \dots, X_n)}{\partial p} \Big|_{p=\hat{p}_{ML}} &= \frac{\sum_{i=1}^n X_i}{\hat{p}_{ML}} - \frac{n - \sum_{i=1}^n X_i}{n - \hat{p}_{ML}} = 0 \\ \Rightarrow \hat{p}_{ML} &= \frac{\sum_{i=1}^n X_i}{n} = \bar{X} \end{aligned}$$

- (b). Suppose $\sum_{i=1}^n X_i = 0$, then

$$l(p \mid X_1, \dots, X_n) = n \log(1 - p), \quad p \in [0, 1] \Rightarrow \hat{p}_{ML} = 0$$

(c). Suppose $\sum_{i=1}^n X_i = n$, then

$$l(p \mid X_1, \dots, X_n) = n \log p, \quad p \in (0, 1] \Rightarrow \hat{p}_{ML} = 1$$

All in all,

$$\hat{p}_{ML} = \bar{X}$$

Remark $\hat{p}_{ML} = \bar{X} = \hat{p}_{MM}$ is the 'best' estimator of p .

2. Suppose $X_i \sim \text{i.i.d. } U[0, \theta]$ where $\theta > 0$ is unknown. The marginal pdf is

$$f(x \mid \theta) = \begin{cases} \frac{1}{\theta}, & x \in [0, \theta] \\ 0, & \text{otherwise} \end{cases} = \frac{1}{\theta} \mathbf{1}_{\{x \in [0, \theta]\}}$$

and the likelihood function is

$$\begin{aligned} L(\theta \mid X_1, \dots, X_n) &= \prod_{i=1}^n \left\{ \frac{1}{\theta} \mathbf{1}_{\{x \in [0, \theta]\}} \right\} = \begin{cases} \frac{1}{\theta^n}, & \theta \geq X_{(n)} \\ 0, & \text{otherwise} \end{cases} \\ &\Rightarrow \hat{\theta}_{ML} = X_{(n)} \end{aligned}$$

Remark $\hat{\theta}_{ML} = X_{(n)} \neq 2\bar{X} = \hat{\theta}_{MM}$; $\hat{\theta}_{ML} < X_i$ can't occur, which is good news; $\hat{\theta}_{ML} \leq \theta$ (low) must occur, which is bad news.

3. Suppose $X_i \sim \text{i.i.d. } N(\mu, \sigma^2)$ where $\mu \in \mathbb{R}$ and $\sigma^2 > 0$ are unknown. Then,

$$\hat{\mu}_{ML} = \hat{\mu}_{MM} = \bar{X}, \quad \hat{\sigma}_{ML}^2 = \hat{\sigma}_{MM}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

2.3.3 Comparing Estimators: Mean Squared Error

General Approach

- o Statistical Decision Theory

Leading Special Case: Mean Squared Error.

Definition 2.10 (Mean Squared Error)

The **mean squared error** (MSE) of one estimator $\hat{\theta}$ of θ is defined as

$$\text{MSE}_\theta(\hat{\theta}) = \mathbb{E}_\theta[(\hat{\theta} - \theta)^2], \quad \theta \in \Theta \subseteq \mathbb{R}$$

Definition 2.11 (Bias)

The **bias** of $\hat{\theta}$ is (the function of θ) given by

$$\text{Bias}_\theta(\hat{\theta}) = \mathbb{E}_\theta(\hat{\theta}) - \theta, \quad \theta \in \Theta$$

$\hat{\theta}$ is **unbiased** iff $\text{Bias}_\theta(\hat{\theta}) = 0 \ (\forall \theta \in \Theta)$

Decomposition:

$$\text{MSE}_\theta(\hat{\theta}) = \text{Bias}_\theta(\hat{\theta})^2 + \text{Var}_\theta(\hat{\theta})$$

which is given by $\mathbb{E}[X^2] = \mathbb{E}[X]^2 + \text{Var}(X)$. Hence, if $\hat{\theta}$ is unbiased ($\text{Bias}_\theta(\hat{\theta}) = 0$), $\text{MSE}_\theta(\hat{\theta}) = \text{Var}_\theta(\hat{\theta})$.

Definition 2.12 (Uniform Minimum Variance Unbiased (UMVU))

An unbiased estimator $\hat{\theta}$ is a **uniform minimum variance unbiased (UMVU)** estimator (of θ) iff

$$\text{MSE}_\theta(\hat{\theta}) = \text{Var}_\theta(\hat{\theta}) \leq \text{Var}_\theta(\tilde{\theta}) = \text{MSE}_\theta(\tilde{\theta})$$

whenever $\tilde{\theta}$ is an unbiased estimator of θ .



Remark UMVU estimators often exist; UMVU estimators are based on sufficient statistics.

2.3.4 Sufficient Statistic**Definition 2.13 (Sufficient Statistic)**

A statistic $T = T(X_1, \dots, X_n)$ is **sufficient** iff the conditional distribution of (X_1, \dots, X_n) given T , $(X_1, \dots, X_n)|T$, doesn't depend on θ .

$$f_X(x | T(X_1, \dots, X_n) = t; \theta) = f_X(x | T(X_1, \dots, X_n) = t), \forall x$$

That is, the mutual information between θ and $T(X_1, \dots, X_n)$ equals the mutual information between θ and $\{X_1, \dots, X_n\}$,

$$I(\theta; T(X_1, \dots, X_n)) = I(\theta; \{X_1, \dots, X_n\})$$

**Theorem 2.4 (Rao-Blackwell Theorem)**

Suppose $\tilde{\theta}$ is an unbiased estimator of θ and suppose T is sufficient (for θ). Then,

- (a). $\hat{\theta} = \mathbb{E}[\tilde{\theta}|T]$ is an unbiased estimator of θ .
- (b). $\text{Var}_\theta(\hat{\theta}) \leq \text{Var}_\theta(\tilde{\theta}), \forall \theta \in \Theta$.

**Proof 2.1**

- (a). Estimator: $\hat{\theta} = \mathbb{E}[\tilde{\theta} | T]$ doesn't depend on θ because T is sufficient. By the Law of Iterative Expectation, we have

$$\mathbb{E}_\theta(\hat{\theta}) = \mathbb{E}_\theta[\mathbb{E}[\tilde{\theta} | T]] = \mathbb{E}_\theta[\tilde{\theta}] = \theta$$

- (b). Variance Reduction: By the Law of Total Variance

$$\text{Var}(\hat{\theta}) = \text{Var}_\theta[\mathbb{E}[\tilde{\theta} | T]] \leq \text{Var}_\theta(\tilde{\theta}), \forall \theta \in \Theta$$

with strict inequality unless $\text{Var}(\tilde{\theta}|T) = 0$ (which also makes $\hat{\theta} = \tilde{\theta}$).

Finding sufficient statistics

- o Apply "definition";
- o Apply factorization criterion.

Proposition 2.3 (Fisher-Neyman Factorization Criterion)

A statistic $T = T(X_1, \dots, X_n)$ is sufficient if and only if $\exists g(\cdot | \cdot)$ and $h(\cdot)$ such that

$$\begin{aligned} f_X((X_1, \dots, X_n) | \theta) &= \prod_{i=1}^n f(X_i | \theta) \\ &= g[T(X_1, \dots, X_n) | \theta] h(X_1, \dots, X_n) \end{aligned}$$



Example 2.6

1. Suppose $\{X_i\}_{i=1}^n$ be a random sample from $Poisson(\theta)$. Then, show $T(X_1, \dots, X_n) = \sum_{i=1}^n X_i$ is a sufficient statistic.

(a). **Prove by Definition:** The sum of independent Poisson random variables are Poisson random variable, so we have $T = \sum_{i=1}^n X_i \sim Pois(n\theta)$. Then the conditional distribution of X_1, \dots, X_n given T is

$$f(X_1, \dots, X_n | T) = \frac{\prod_{i=1}^n \frac{\theta^{X_i} e^{-\theta}}{X_i!}}{\frac{(n\theta)^T e^{-n\theta}}{T!}} = \frac{T!}{n^T \prod_{i=1}^n X_i!}$$

which is independent of θ . So, $T(X_1, \dots, X_n)$ is sufficient statistic by definition.

(b). **Prove by Factorization Theorem:**

$$\prod_{i=1}^n f(X_i | \theta) = \prod_{i=1}^n \frac{\theta^{X_i} e^{-\theta}}{X_i!} = \frac{\theta^{T(X_1, \dots, X_n)} e^{-n\theta}}{\prod_{i=1}^n X_i!} = g(T(X_1, \dots, X_n) | \theta) h(X_1, \dots, X_n)$$

where $g(T(X_1, \dots, X_n) | \theta) = \theta^{T(X_1, \dots, X_n)} e^{-n\theta}$ and $h(X_1, \dots, X_n) = \frac{1}{\prod_{i=1}^n X_i!}$. Hence, $T(X_1, \dots, X_n)$ is sufficient statistic by Fisher-Neyman Factorization Criterion.

(c). **Prove by Exponential Family:**

$$f(X | \theta) = \frac{\theta^X e^{-\theta}}{X!} = \frac{e^{-\theta + X \ln \theta}}{X!}$$

Hence, the distribution is a member of the exponential family, where $c(\theta) = 1$, $h(X) = \frac{1}{X!}$, $w_1(\theta) = -\theta$, $w_2(\theta) = \ln \theta$, $t_1(X) = 1$, $t_2(X) = X$. By theorem 2.7, $\sum_{i=1}^n X_i$ is sufficient because $\{w_1(\theta) = -\theta, w_2(\theta) = \ln \theta\}$ is non-empty.

2. Suppose $X_i \sim$ i.i.d. $U[0, \theta]$ where $\theta > 0$ is unknown. The marginal pdf is

$$f(x | \theta) = \begin{cases} \frac{1}{\theta}, & x \in [0, \theta] \\ 0, & \text{otherwise} \end{cases} = \frac{1}{\theta} \mathbf{1}_{\{x \in [0, \theta]\}}$$

Factorization:

$$\prod_{i=1}^n f(X_i | \theta) = \underbrace{\frac{1}{\theta^n} \mathbf{1}_{\{X_{(n)} \leq \theta\}}}_{g(X_{(n)}|\theta)} \underbrace{\mathbf{1}_{\{X_{(1)} \geq 0\}}}_{h(X_1, \dots, X_n)}$$

Hence, we have shown that $X_{(n)}$ is sufficient $\Rightarrow \hat{\theta}_{MM} = 2\bar{X}$ cannot be UMVU and $\hat{\theta}_{RB} = \mathbb{E}[\hat{\theta}_{MM}|X_{(n)}]$ is better.

2.3.5 Minimal Sufficient Statistic

Definition 2.14 (Minimal Sufficient Statistic)

A sufficient statistic $T(X_1, \dots, X_n)$ is called a **minimal sufficient statistic** if, for any other sufficient statistic $T'(X_1, \dots, X_n)$, $T(X_1, \dots, X_n)$ is a function of $T'(X_1, \dots, X_n)$.



Theorem 2.5 (Theorem to Check Minimal Sufficient Statistic)

Let $f(\vec{X})$ be the pmf or pdf of a sample \vec{X} . Suppose there exists a function $T(\vec{X})$ such that,

"for every sample points \vec{X} and \vec{Y} , the ratio $\frac{f(\vec{X}|\theta)}{f(\vec{Y}|\theta)}$ is constant for any θ if and only if $T(\vec{X}) = T(\vec{Y})$ ".

Then $T(\vec{X})$ is a **minimal sufficient statistic** for θ .



Example 2.7 Let $X_1, \dots, X_n \sim \text{i.i.d. } U[\theta - \frac{1}{2}, \theta + \frac{1}{2}]$, with $\theta \in \mathbb{R}$ unknown.

By $f(X | \theta) = \mathbf{1}_{\{X \in [\theta - \frac{1}{2}, \theta + \frac{1}{2}]\}}$, we have

$$\prod_{i=1}^n f(X_i | \theta) = \underbrace{\mathbf{1}_{\{X_{(1)} \geq \theta - \frac{1}{2}\}} \mathbf{1}_{\{X_{(n)} \leq \theta + \frac{1}{2}\}}}_{g[T(X_1, \dots, X_n)|\theta]} \underbrace{\mathbf{1}}_{h(X_1, \dots, X_n)}$$

By the Fisher-Neyman Factorization Criterion, $T(X_1, \dots, X_n) = \{X_{(1)}, X_{(n)}\}$ is a sufficient statistic.

We can prove $T(X_1, \dots, X_n) = \{X_{(1)}, X_{(n)}\}$ is a minimal sufficient statistic by proving "for every sample points (X_1, \dots, X_n) and (Y_1, \dots, Y_n) , $\frac{f(X_1, \dots, X_n|\theta)}{f(Y_1, \dots, Y_n|\theta)}$ is constant as a function of θ if and only if $T(X_1, \dots, X_n) = T(Y_1, \dots, Y_n)$."

$$\frac{f(X_1, \dots, X_n | \theta)}{f(Y_1, \dots, Y_n | \theta)} = \frac{\mathbf{1}_{\{X_{(1)} \geq \theta - \frac{1}{2}\}} \mathbf{1}_{\{X_{(n)} \leq \theta + \frac{1}{2}\}}}{\mathbf{1}_{\{Y_{(1)} \geq \theta - \frac{1}{2}\}} \mathbf{1}_{\{Y_{(n)} \leq \theta + \frac{1}{2}\}}}$$

Hence, for every sample points (X_1, \dots, X_n) and (Y_1, \dots, Y_n) , $\frac{f(X_1, \dots, X_n|\theta)}{f(Y_1, \dots, Y_n|\theta)}$ is constant for all θ if and only if $X_{(1)} = Y_{(1)}$ and $X_{(n)} = Y_{(n)}$. That is, $T(X_1, \dots, X_n) = T(Y_1, \dots, Y_n)$. Hence, $T(X_1, \dots, X_n) = \{X_{(1)}, X_{(n)}\}$ is a **minimal sufficient statistic**.

Consider $g(T) = X_{(n)} - X_{(1)} - \frac{n-1}{n+1}$, it has $\mathbb{E}[g(T)] = 0$ but $P_\theta[g(T) = 0] < 1$. Hence, T is not a complete statistic by definition.

2.3.6 Complete Statistic

Suppose T is sufficient and then $\hat{\theta} = \hat{\theta}(T)$ is unbiased. Under what conditions (on T) is $\hat{\theta}$ UMVU?

Answers: If "only one" estimator based on T is unbiased. (T is complete.)

Definition 2.15 (Complete Statistic)

A statistic T is **complete** if and only if

$$P_{\theta}[g(T) = 0] = 1, \forall \theta \in \Theta$$

whenever $g(\cdot)$ is such that

$$\mathbb{E}_{\theta}[g(T)] = 0, \forall \theta \in \Theta$$

(whenever the mean is zero, it can only equal to zero). 

Recall: A matrix $A_{m \times k}$ has rank k iff $Ax = 0 \Rightarrow x = 0$.

Theorem 2.6 (Lehmann-Scheffé Theorem)

If T is complete and if $\hat{\theta} = \hat{\theta}(T)$ and $\tilde{\theta} = \tilde{\theta}(T)$ are unbiased, then

$$\mathbb{E}_{\theta}[\hat{\theta} - \tilde{\theta}] = 0 \Rightarrow P(\hat{\theta} - \tilde{\theta} = 0) = P(\hat{\theta} = \tilde{\theta}) = 1$$


Implication:

Corollary 2.1 (Unbiased $\hat{\theta}(T)$ with sufficient and complete T)

If T is sufficient and complete and if $\hat{\theta} = \hat{\theta}(T)$ is unbiased, then $\hat{\theta}$ is UMVU (let $\tilde{\theta}$ be an UMVU). 

Example 2.8 Suppose $X_i \sim \text{i.i.d. } U[0, \theta]$ where $\theta > 0$ is unknown.

Facts:

- $X_{(n)}$ is sufficient and complete \Rightarrow Any unbiased estimator given $X_{(n)}$ is UMVU, e.g. $\hat{\theta}_{RB} = \mathbb{E}[\hat{\theta}_{MM}|X_{(n)}]$;
- $\mathbb{E}_{\theta}(X_{(n)}) = \frac{n}{n+1}\theta \Rightarrow$ unbiased $\frac{n+1}{n}X_{(n)}$ is UMVU ($= \hat{\theta}_{RB}$).

Remark The cdf of $X_{(n)}$ is

$$F_{X_{(n)}}(x | \theta) = F(x | \theta)^n = \begin{cases} 0, & \text{if } x < 0 \\ \left(\frac{x}{\theta}\right)^n & \text{if } 0 \leq x \leq \theta \\ 1, & \text{if } x > \theta \end{cases}$$

so $X_{(n)}$ is continuous with pdf

$$f_{X_{(n)}}(x | \theta) = \begin{cases} \frac{n}{\theta^n} x^{n-1} & \text{if } x \in [0, \theta] \\ 0, & \text{otherwise} \end{cases}$$

Hence, $\mathbb{E}_{\theta} X_{(n)} = \int_0^{\theta} \frac{n}{\theta^n} x^{n-1} x dx = \frac{n}{n+1}\theta$.

Verifying Completeness

- Apply definition:
 - Example: $\sum_{i=1}^n X_i$ is complete when $X_i \sim i.i.d. \text{Ber}(p)$ - compute rank of the matrix to check completeness
- Show that $\{f(\cdot|\theta) : \theta \in \Theta\}$ is on exponential family and apply theorem 2.7.

Theorem 2.7 (Sufficient and Complete Statistic for Exponential Family)

If the distribution is a member of the exponential family, that is,

$$f(x|\theta) = c(\theta)h(x)\exp\left\{\sum_{j=1}^k w_j(\theta)t_j(x)\right\}$$

then

$$T = \left(\sum_{i=1}^n t_1(x_i), \dots, \sum_{i=1}^n t_k(x_i) \right)$$

is sufficient and complete if $\{w_1(\theta), \dots, w_k(\theta)\} : \theta \in \Theta$ contains an open set.



Example 2.9 Suppose $X \sim \mathcal{N}(\mu, \sigma^2)$ for some $\mu \in \mathbb{R}$ and some $\sigma^2 > 0$. Then, $\theta = (\mu, \sigma^2)$ and $\Theta = \mathbb{R} \times \mathbb{R}_{++}$.

The pdf can be written as

$$f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{\mu^2}{2\sigma^2}} e^{\frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2}$$

We can have $h(x) = 1, c(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{\mu^2}{2\sigma^2}}, t_1(x) = x, w_1(\mu, \sigma^2) = \frac{\mu}{\sigma^2}, t_2(x) = x^2, w_2(\mu, \sigma^2) = -\frac{1}{2\sigma^2}$.

That is, $T = (\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2)$ is sufficient and complete.

And $(\bar{X}, S^2) = \left(\frac{1}{n} \sum_{i=1}^n X_i, \frac{1}{n-1} \sum_{i=1}^n \left[X_i^2 - \frac{(\sum_{i=1}^n X_i^2)^2}{n} \right] \right)$ is UMVU estimator of (μ, σ^2) .

2.3.7 Cramér-Rao Lower Bound

The score function is the derivative of the log likelihood function with respect to θ .

Definition 2.16 (Score Function)

The **score function** is

$$u(\theta, \vec{X}) = \frac{\partial}{\partial \theta} \log f_{\vec{X}}(\vec{X} | \theta) = \frac{1}{f_{\vec{X}}(\vec{X} | \theta)} \frac{\partial}{\partial \theta} f_{\vec{X}}(\vec{X} | \theta)$$

where $f_{\vec{X}}(\vec{X} | \theta) = L(\theta | X_1, \dots, X_n) = \prod_{i=1}^n f_{X_i}(X_i | \theta)$.



Definition 2.17 ("Regularity" Condition)

The regularity conditions are as follows:

1. The partial derivative of $f_{\vec{X}}(\vec{X} | \theta)$ with respect to θ exists almost everywhere. (It can fail to exist on a null set, as long as this set does not depend on θ .)

2. The integral of $f_{\vec{X}}(\vec{X} | \theta)$ can be differentiated under the integral sign with respect to θ .
3. The support of $f_{\vec{X}}(\vec{X} | \theta)$ does not depend on θ .



Lemma 2.1 (Mean of Score Function is Zero)

Under "Regularity" condition and X are continuous, the mean of score function is zero:

$$\begin{aligned}\mathbb{E}_\theta[u(\theta, \vec{X})] &= \int_{\vec{X}} \left[\frac{\partial}{\partial \theta} \log f_{\vec{X}}(\vec{X} | \theta) \right] f_{\vec{X}}(\vec{X} | \theta) d\vec{X} \\ &= \int_{\vec{X}} \left[\frac{\partial}{\partial \theta} f_{\vec{X}}(\vec{X} | \theta) \right] d\vec{X} \\ (*) \quad &= \frac{\partial}{\partial \theta} \underbrace{\int_{\vec{X}} f_{\vec{X}}(\vec{X} | \theta) d\vec{X}}_{=1} = 0\end{aligned}$$

(*): Moving the derivative outside the integral can be done as long as the limits of integration are fixed, i.e. they do not depend on θ .



Definition 2.18 (Fisher Information)

Under "regularity" conditions, the **Fisher information** is defined to be the variance of the score function.

$$\mathcal{I}(\theta) = \text{Var}_\theta(u(\theta, \vec{X})) = \mathbb{E}_\theta[u(\theta, \vec{X})u(\theta, \vec{X})^T]$$



Lemma 2.2 (Fisher Information)

For $\theta \in \Theta \subseteq \mathbb{R}$, the Fisher information can be written as

$$\begin{aligned}\mathcal{I}(\theta) &= \text{Var}_\theta \left[\frac{\partial}{\partial \theta} \log f_{\vec{X}}(\vec{X} | \theta) \right] \\ &= \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial \theta} \log f_{\vec{X}}(\vec{X} | \theta) \right)^2 \right] = n \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial \theta} \log f_X(X_i | \theta) \right)^2 \right] \\ &= -\mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \log f_{\vec{X}}(\vec{X} | \theta) \right] = -n \mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \log f_X(X_i | \theta) \right]\end{aligned}$$



Proposition 2.4 (Cramér-Rao Lower Bound)

Under "regularity" conditions, for every estimator $\hat{\theta}$

$$\text{Var}_\theta[\hat{\theta}(\vec{X})] \geq \frac{\left(\frac{d}{d\theta} \mathbb{E}_\theta[\hat{\theta}(\vec{X})] \right)^2}{\mathcal{I}(\theta)} \equiv \text{CRLB}(\theta)$$

Specifically, if the estimator $\hat{\theta}$ is unbiased,

$$\text{CRLB}(\theta) = \mathcal{I}(\theta)^{-1}$$



Remark $\mathcal{I}(\theta)$ is called the **Fisher Information**; "Regularity" conditions are satisfied by "smooth" exponential families; Proof uses Cauchy-Schwarz inequality.

3 Possibilities

(1). CR inequality is applicable and attainable:

- (a). Estimating p when $X \sim$ i.i.d. $\text{Ber}(p)$;
- (b). Estimating μ when $X \sim$ i.i.d. $N(\mu, \sigma^2)$.

(2). CR inequality is applicable, but not attainable:

- (a). Estimating σ^2 when $X \sim$ i.i.d. $N(\mu, \sigma^2)$: $\text{Var}(S^2) = \frac{2\sigma^4}{n-1} > \frac{2\sigma^4}{n} = \mathcal{I}(\theta)^{-1}$ (CR bound).

(3). CR inequality is not applicable:

- (a). Estimating θ when $X \sim$ i.i.d. $U[0, \theta]$: CR bound $\mathcal{I}(\theta)^{-1} = \frac{\theta^2}{n}$ and $\text{Var}(\hat{\theta}_{UMVU}) = \frac{\theta^2}{n(n+2)}$

Theorem 2.8 (MLE Covariance $\xrightarrow{n \rightarrow \infty}$ Cramér-Rao Lower Bound)

Suppose the sample $\{X_i\}_{i=1}^n$ is i.i.d. The Maximum likelihood estimator (MLE) $\hat{\theta} = \arg \max_{\theta} L(\theta \mid X_1, \dots, X_n)$, under "regularity" conditions, as $n \rightarrow \infty$

$$\sqrt{n} (\hat{\theta} - \theta) \xrightarrow{D} N(0, \mathcal{I}(\theta)^{-1})$$



Proposition 2.5 (Approximation of MLE Covariance Matrix)

When the sample x is made up of i.i.d. observations, the covariance matrix of the maximum likelihood estimator $\hat{\theta}$ is approximately equal to the inverse of the information matrix.

$$\text{Cov}(\hat{\theta}) \approx (\mathcal{I}(\theta))^{-1}$$



Hence, the covariance matrix can be estimated as $(\mathcal{I}(\hat{\theta}))^{-1}$. Similarly, SE is estimated by $\sqrt{(\mathcal{I}(\hat{\theta}))^{-1}}$.

2.4 Hypothesis Testing

X_1, \dots, X_n is a random sample from a discrete/continuous distribution with pmf/pdf $f(\cdot \mid \theta)$, where $\theta \in \Theta$ is unknown.

Ingredients of Hypothesis Test

- (1). Formulation of Testing Problem:
 - Partitioning of Θ into two disjoint subsets Θ_0 and Θ_1 .
- (2). Testing Procedure:
 - Rule for choosing the two subsets specified in (1).

2.4.1 Formulation of Testing Problem

Formulating a Testing Procedure

- Terminology:

Definition 2.19 (Hypothesis)

- (a). A hypothesis is a statement about θ ;
- (b). Null hypothesis: $H_0 : \theta \in \Theta_0$;
- (c). Alternative hypothesis: $H_1 : \theta \in \Theta \setminus \Theta_0$;
- (d). Maintained hypothesis: $\theta \in \Theta$ (always correct).
- (e). Typical Formulation:

$$H_0 : \theta \in \Theta_0 \text{ vs. } H_1 : \theta \in \Theta_1$$



Example 2.10 Suppose $X \sim \text{i.i.d. } N(\mu, 1)$, where $\mu \geq 0$ is unknown.

Objective: Determine whether $\mu = 0$.

Two possible formulation: $H_0 : \mu = 0$ vs. $H_1 : \mu > 0$ (or vice versa).

- Testing Procedure:

Consider the problem of testing $H_0 : \theta \in \Theta_0$ vs. $H_1 : \theta \in \Theta_1$.

Definition 2.20 (Testing Procedure with Critical Region)

A testing procedure is a (data-based) rule for choosing between H_0 and H_1 .

The rule:

”Reject H_0 iff $(X_1, \dots, X_n) \in C$ ” (for some $C \in \mathbb{R}^n$)

is a testing procedure with critical region C .



Example 2.11 Suppose $X \sim \text{i.i.d. } N(\mu, 1)$, where $\mu \geq 0$ is unknown. The decision rule ”Reject H_0 iff $\frac{\sum_{i=1}^n X_i}{n} = \bar{X} \geq \frac{1.645}{\sqrt{n}}$ ”, where the critical region is $C = \{(X_1, \dots, X_n) : \frac{\sum_{i=1}^n X_i}{n} \geq \frac{1.645}{\sqrt{n}}\}$

Proposition 2.6 (Critical Region \Leftrightarrow Test Statistic and Critical Value)

Any set $C \in \mathbb{R}^n$ can be written as

$$C = \{(X_1, \dots, X_n) : T(X_1, \dots, X_n) > c\}$$

for some $T : \mathbb{R}^n \rightarrow \mathbb{R}$ and some $c \in \mathbb{R}$.

**Definition 2.21 (Test Statistic and Critical Value)**

$T(X_1, \dots, X_n)$ is called a test statistic and c is called the critical value (of the test).



2.4.2 Errors, Power Function, and Agenda

Agenda

1. Choosing critical value (given test statistic).
2. Choosing test statistic.

Definition 2.22 (Type I and Type II Errors)

Decision vs. Truth	H_0 (True)	H_1 (False)
H_0 (Fail to Reject)		Type II Error
H_1 (Reject)		Type I Error

where

1. Type I Error: mistaken rejection of a null hypothesis that is actually true;
2. Type II Error: failure to reject a null hypothesis that is actually false.



There is a trade-off between Type I and Type II errors. The general approach is *statistical decision theory*.

Example 2.12 Heading Special Case: Making $P_\theta[\text{Type I Error}]$ "small".

Definition 2.23 (Power Function)

The **power function** of a test unit critical region $C \subseteq \mathbb{R}^n$ is the function $\beta : \Theta \rightarrow [0, 1]$ given by

$$\begin{aligned}\beta(\theta) &= P_\theta[\text{Reject } H_0] \\ &= P_\theta[(X_1, \dots, X_n)' \in C] \\ (\text{equivalently}) &= P_\theta[T(X_1, \dots, X_n) > c]\end{aligned}$$

for corresponding statistic T and critical value c .



- For $\theta \in \Theta_0$: $P_\theta[\text{Type I Error}] = P_\theta[\text{Reject } H_0] = \beta(\theta)$;
- For $\theta \in \Theta_1$: $P_\theta[\text{Type II Error}] = 1 - P_\theta[\text{Reject } H_0] = 1 - \beta(\theta)$;
- Hence, the ideal power function is $\beta(\theta) = \begin{cases} 1, & \theta \in \Theta_1; \\ 0, & \theta \in \Theta_0 \end{cases}$;
- "Good" Power Function: $\beta(\theta)$ is "low" ("high") when $\theta \in \Theta_0$ ($\theta \in \Theta_1$).

Standard:

- (1). Given $T(\cdot)$, choose critical value c such that $\beta(\theta) = P_\theta[T(X_1, \dots, X_n) > c] \leq 5\%$ when $\theta \in \Theta_0$ (i.e., $\sup_{\theta \in \Theta_0} \beta(\theta) \leq 5\%$);
- (2). Choose test statistic such that $\beta(\theta) = P_\theta[T(X_1, \dots, X_n) > c(T)]$ is "large" for $\theta \in \Theta_1$. (Main Tool: Neyman-Pearson Lemma).

2.4.3 Choice of Critical Value

Given $T(\cdot)$, choose critical value c such that $\beta(\theta) = P_\theta[T(X_1, \dots, X_n) > c] \leq 5\%$ when $\theta \in \Theta_0$ (i.e., $\sup_{\theta \in \Theta_0} \beta(\theta) \leq 5\%$).

Definition 2.24 (Test Size and Level α)

The **size** of a test (with power function β) is $\sup_{\theta \in \Theta_0} \beta(\theta)$.

A test is of **level α** ($\in [0, 1]$) if and only if its size is $\leq \alpha$. (Standard choice $\alpha = 0.05$). 

Example 2.13 Suppose $X \sim \text{i.i.d. } N(\mu, 1)$, where $\mu \geq 0$ is unknown.

Consider the decision rule "Reject H_0 iff $\frac{\sum_{i=1}^n X_i}{n} = \bar{X} \geq \frac{1.645}{\sqrt{n}}$ ". The power function is $\beta(\mu) = P_\mu[\text{Reject } H_0] = P_\mu(\bar{X} \geq \frac{1.645}{\sqrt{n}})$

Recall: $\bar{X} \sim \mathcal{N}(\mu, \frac{1}{n}) \Rightarrow \sqrt{n}(\bar{X} - \mu) \sim \mathcal{N}(0, 1)$.

$$\begin{aligned}\beta(\mu) &= P_\mu[\text{Reject } H_0] = P_\mu(\bar{X} \geq \frac{1.645}{\sqrt{n}}) \\ &= P_\mu(\sqrt{n}(\bar{X} - \mu) \geq 1.645 - \sqrt{n}\mu) \\ &= 1 - \Phi(1.645 - \sqrt{n}\mu)\end{aligned}$$

where Φ is the standard normal cdf.

Size = $\beta(0) = 1 - \Phi(1.645) \approx 0.05$.

2.4.4 Choice of Test Statistic: Uniformly Most Powerful (UMP) Level α Test

Choose test statistic such that $\beta(\theta) = P_\theta[T(X_1, \dots, X_n) > c(T)]$ is "large" for $\theta \in \Theta_1$. (Main Tool: Neyman-Pearson Lemma).

Definition 2.25 (Uniformly Most Powerful (UMP) Level α Test)

A test with level α and power function β is a **uniformly most powerful (UMP) level α test** iff

$$\beta(\theta) \geq \tilde{\beta}(\theta), \forall \theta \in \Theta_1$$

where $\tilde{\beta}$ is the power function of some (other) level α test. 

Consider the problem of testing $H_0 : \theta = \theta_0 \in \mathbb{R}$

- UMP level α test always \exists if $H_1 : \theta = \theta_1$ (Proven by Neyman-Pearson Lemma);
- UMP level α test often \exists if $H_1 : \theta > \theta_0$ or $H_1 : \theta < \theta_0$ (Proven by Karlin-Rubin Theorem);
- UMP level α test often \nexists if $H_1 : \theta \neq \theta_0$; UMP "unbiased" level α test often \exists .

Theorem 2.9 (Neyman-Pearson Lemma)

Consider the problem of testing,

$$H_0 : \theta = \theta_0 \text{ vs. } H_1 : \theta = \theta_1$$

For any $k \geq 0$, the test which

$$\text{Rejects } H_0 \text{ iff } L(\theta_1 | X_1, \dots, X_n) \geq kL(\theta_0 | X_1, \dots, X_n)$$

is a UMP level α test, where

$$\alpha = P_{\theta_0}[L(\theta_1 | X_1, \dots, X_n) \geq k L(\theta_0 | X_1, \dots, X_n)]$$

and where $L(\theta | X_1, \dots, X_n) = \prod_{i=1}^n f(X_i | \theta)$.



Remark

- UMP level α test exists if $\alpha \in \{P_{\theta_0}[L(\theta_1 | X_1, \dots, X_n) \geq k L(\theta_0 | X_1, \dots, X_n)] : k \geq 0\}$.
- The Neyman-Pearson Lemma rejects the H_0 iff

$$L(\theta_1 | X_1, \dots, X_n) \geq k L(\theta_0 | X_1, \dots, X_n) \Leftrightarrow \frac{L(\theta_1 | X_1, \dots, X_n)}{L(\theta_0 | X_1, \dots, X_n)} \geq k$$

$$(L(\theta_0 | X_1, \dots, X_n) \neq 0)$$

- Hence, it is called **"Likelihood Ratio" test.**
- Converse: Any UMP level α test is of "NP type."

Example of Using NP Lemma

Example 2.14 Suppose $X \sim \text{i.i.d. } N(\mu, 1)$, where $\mu \geq 0$ is unknown.

Let $\mu_1 = 0$ be given and consider the problem of testing

$$H_0 : \mu = 0 \text{ vs. } H_1 : \mu = \mu_1 > 0$$

We have $L(\mu | X_1, \dots, X_n) = \prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{(X_i - \mu)^2}{2}} \right) = (2\pi)^{-\frac{n}{2}} e^{-\frac{1}{2} \sum_{i=1}^n X_i^2} e^{\mu \sum_{i=1}^n X_i} e^{-\frac{n\mu^2}{2}}$. Then,

$$\frac{L(\mu = \mu_1 | X_1, \dots, X_n)}{L(\mu = 0 | X_1, \dots, X_n)} = e^{\mu_1 \sum_{i=1}^n X_i} e^{-\frac{n\mu_1^2}{2}}$$

Decision Rule: Reject H_0 iff

$$\begin{aligned} \frac{L(\mu = \mu_1 | X_1, \dots, X_n)}{L(\mu = 0 | X_1, \dots, X_n)} &= e^{\mu_1 \sum_{i=1}^n X_i} e^{-\frac{n\mu_1^2}{2}} \geq k \\ &\Leftrightarrow -\frac{n\mu_1^2}{2} + \mu_1 \sum_{i=1}^n X_i \geq \log k \\ &\Leftrightarrow \bar{X} \geq \frac{\log k}{n\mu_1} + \frac{\mu_1}{2} \end{aligned}$$

The NP test reject for large values of \bar{X} .

Optimality Theorem for One-sided Testing Problem

Consider

$$H_0 : \mu = \mu_0 \text{ vs. } H_1 : \mu > \mu_0$$

For any $\theta_1 > \theta_0$, use NP Lemma to find optimal test of $H_0 : \mu = \theta_0$ vs. $H_1 : \mu = \theta_1$.

- If the NP tests coincide, then the test is the UMP test of $H_0 : \mu = \mu_0$ vs. $H_1 : \mu > \mu_0$;
- Otherwise, \nexists UMP (level α) test of the $H_0 : \mu = \mu_0$ vs. $H_1 : \mu > \mu_0$.

Implications: (The previous $N(\mu, 1)$ example)

- (i). The UMP 5% test of $H_0 : \mu = 0$ vs. $H_1 : \mu > 0$ rejects H_0 iff $\bar{X} > \frac{1.645}{\sqrt{n}}$.
- (ii). The UMP 5% test of $H_0 : \mu = 0$ vs. $H_1 : \mu < 0$ rejects H_0 iff $-\bar{X} > \frac{1.645}{\sqrt{n}}$.
- (iii). \nexists UMP 5% test of $H_0 : \mu = 0$ vs. $H_1 : \mu \neq 0$.

Definition 2.26 (Unbiased Test)

A test of

$$H_0 : \theta \in \Theta_0 \text{ vs. } H_1 : \theta \in \Theta_1$$

is **unbiased** iff its power function $\beta(\cdot)$ satisfies $\sup_{\theta \in \Theta_0} \beta(\theta) \leq \inf_{\theta \in \Theta_1} \beta(\theta)$



Claim 2.1

The UMP unbiased 5% test of $H_0 : \mu = 0$ vs. $H_1 : \mu \neq 0$: Rejects H_0 iff $|\bar{X}| > \frac{1.96}{\sqrt{n}}$.



Corollary 2.2

Suppose $X_i \sim \text{i.i.d. } N(\mu, \sigma^2)$, where σ^2 is known. Then, the UMP unbiased 5% test of the $H_0 : \mu = \mu_0$ vs. $H_1 : \mu \neq \mu_0$: Rejects H_0 if $|\frac{\bar{X} - \mu_0}{\sigma}| > \frac{1.96}{\sqrt{n}}$.



Claim 2.2

"In general", "Natural" test statistics are (approximately) optimal and critical values can be find.



2.4.5 Generalized Neyman-Pearson Lemma

NP Lemma: $\max \beta(\theta_1)$ s.t. $\beta(\theta_0) \leq \alpha$;

Generalized NP Lemma: How to optimize a function with infinity constraints.

Observation: If β is differentiable, then an unbiased test of the $H_0 : \theta = \theta_0$ vs. $H_1 : \theta \neq \theta_0$ satisfies $\beta'(\theta_0) = 0$

Theorem 2.10 (Generalized Neyman-Pearson Lemma)



2.5 Trinity of Classical Tests

- Likelihood Ratio Test
- Lagrangian Multiplier Test (Score Test)
- Wald Test

Properties: Deliver optimal test in motivating example; closely related (and "approximately" optimal) in general.

2.5.1 Test Statistics

Settings: X_1, \dots, X_n is a random sample from a discrete/continuous distribution with pmf/pdf $f(\cdot | \theta)$, where $\theta \in \Theta \subseteq \mathbb{R}$ is unknown.

Testing Problem: $H_0 : \theta = \theta_0$ vs. $H_1 : \theta \neq \theta_0$ for some $\theta_0 \in \Theta$.

Recall the log likelihood function is given by

$$l(\theta | X_1, \dots, X_n) = \sum_{i=1}^n \log f(X_i | \theta)$$

The (sample) score function is

$$u(\theta | X_1, \dots, X_n) = \frac{\partial}{\partial \theta} l(\theta | X_1, \dots, X_n)$$

and the (sample) fisher information is

$$\mathcal{I}(\theta | X_1, \dots, X_n) = -\frac{\partial^2}{\partial \theta^2} l(\theta | X_1, \dots, X_n)$$

- **Likelihood Ratio Test Statistic:**

$$\begin{aligned} T_{LR}(X_1, \dots, X_n) &= 2 \left\{ \max_{\theta \in \Theta} l(\theta | X_1, \dots, X_n) - \max_{\theta \in \Theta_0} l(\theta | X_1, \dots, X_n) \right\} \text{ (general form)} \\ &= 2 \left\{ l(\hat{\theta}_{ML} | X_1, \dots, X_n) - l(\theta_0 | X_1, \dots, X_n) \right\} \\ &= 2 \log \left\{ \frac{l(\hat{\theta}_{ML} | X_1, \dots, X_n)}{l(\theta_0 | X_1, \dots, X_n)} \right\} \end{aligned}$$

Motivation: Neyman-Pearson Lemma (2.9)

- **Lagrangian Multiplier Test Statistic:**

$$T_{LM}(X_1, \dots, X_n) = \frac{\left(\frac{\partial}{\partial \theta} l(\theta_0 | X_1, \dots, X_n) \right)^2}{-\frac{\partial^2}{\partial \theta^2} l(\theta_0 | X_1, \dots, X_n)} = \frac{(u(\theta_0 | X_1, \dots, X_n))^2}{\mathcal{I}(\theta_0 | X_1, \dots, X_n)}$$

Motivation: T_{LM} is approximate to T_{LR} ; No estimation required.

- **Wald Test Statistic:**

$$T_W(X_1, \dots, X_n) = \frac{(\hat{\theta}_{ML} - \theta_0)^2}{\left\{ -\frac{\partial^2}{\partial \theta^2} l(\hat{\theta}_{ML} | X_1, \dots, X_n) \right\}^{-1}} = \frac{(\hat{\theta}_{ML} - \theta_0)^2}{\left(\mathcal{I}(\hat{\theta}_{ML} | X_1, \dots, X_n) \right)^{-1}}$$

Motivation: T_W is approximate to T_{LR} ;

Generalization: Reject the $H_0 : \theta = \theta_0$ if $|\hat{\theta} - \theta_0|$ is "large", when $\hat{\theta}$ is some estimator of θ .

Claim 2.3

In general, for "large" n ,

$$T_{LR} \approx T_{LM} \approx T_W \sim \chi^2(1) = N(0, 1)^2 \text{ under } H_0 : \theta = \theta_0$$

- Approximate 5% critical value is $(1.96)^2 = 3.84$.
- $T_{LR} = T_{LM} = T_W \sim \chi^2(1) = N(0, 1)^2$ under $H_0 : \theta = \theta_0$ when $X_i \sim \text{i.i.d. } N(\mu, 1)$.



2.5.2 Approximation to T_{LR}

In this part as $n \rightarrow \infty$, we use $l(\theta), l'(\theta), l''(\theta)$ to denote $l(\theta \mid X_1, \dots, X_n), l'(\theta \mid X_1, \dots, X_n) \triangleq u(\theta \mid X_1, \dots, X_n), l''(\theta \mid X_1, \dots, X_n) \triangleq -\mathcal{I}(\theta \mid X_1, \dots, X_n)$.

(1). T_{LM} :

Suppose

$$l(\theta) \approx l(\theta_0) + l'(\theta_0)(\theta - \theta_0) + \frac{1}{2}l''(\theta_0)(\theta - \theta_0)^2 \triangleq \tilde{l}(\theta)$$

Then

$$\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmax}} l(\theta) \approx \underset{\theta}{\operatorname{argmax}} \tilde{l}(\theta) = \theta_0 - \frac{l'(\theta_0)}{l''(\theta_0)} \triangleq \tilde{\theta}_{ML}$$

Hence,

$$T_{LR} = 2 \left\{ l(\hat{\theta}_{ML}) - l(\theta_0) \right\} \approx 2 \left\{ \tilde{l}(\tilde{\theta}_{ML}) - \tilde{l}(\theta_0) \right\} = -\frac{l'(\theta_0)^2}{l''(\theta_0)} = T_{LM}$$

(2). T_W :

Suppose

$$l(\theta) \approx l(\hat{\theta}_{ML}) + l'(\hat{\theta}_{ML})(\theta - \hat{\theta}_{ML}) + \frac{1}{2}l''(\hat{\theta}_{ML})(\theta - \hat{\theta}_{ML})^2 \triangleq \hat{l}(\theta)$$

Then,

$$T_{LR} = 2 \left\{ l(\hat{\theta}_{ML}) - l(\theta_0) \right\} \approx 2 \left\{ \tilde{l}(\hat{\theta}_{ML}) - \hat{l}(\theta_0) \right\} = \frac{(\hat{\theta}_{ML} - \theta_0)^2}{(-l''(\hat{\theta}_{ML}))^{-1}} = T_W$$

2.6 Interval Estimation

Definition 2.27

Suppose $\theta \in \mathbb{R}$.

1. An interval estimator of θ is an interval $[L(X_1, \dots, X_n), U(X_1, \dots, X_n)]$, where $L(X_1, \dots, X_n)$ and $U(X_1, \dots, X_n)$ are statistics.
2. The converge probability (of the interval estimator) is the function (of θ) given by

$$P_\theta [L(X_1, \dots, X_n) \leq \theta \leq U(X_1, \dots, X_n)]$$

3. The confidence coefficient is $\inf_\theta P_\theta [L(X_1, \dots, X_n) \leq \theta \leq U(X_1, \dots, X_n)]$



Example 2.15 Suppose $X_i \sim \text{i.i.d. } N(\mu, 1)$, where μ is unknown.

Interval estimator: $\left[\bar{X} - \frac{1.96}{\sqrt{n}}, \bar{X} + \frac{1.96}{\sqrt{n}} \right]$.

Converge probability: $P_\mu \left[\bar{X} - \frac{1.96}{\sqrt{n}} \leq \mu \leq \bar{X} + \frac{1.96}{\sqrt{n}} \right] = P_\mu \left[-1.96 \leq \sqrt{n}(\bar{X} - \mu) \leq 1.96 \right] = \Phi(1.96) - \Phi(-1.96) \approx 0.95$.

Interpretation:

(I). Recall

- (i). $\bar{X} = \hat{\mu}_{MM} = \hat{\mu}_{ML} = \hat{\mu}_{UMVU}$;
- (ii). $\bar{X} \sim \mathcal{N}(\mu, \frac{1}{n}) \Rightarrow \frac{1}{\sqrt{n}} = \sqrt{\text{Var}(\bar{x})}$.

$$\text{Hence, } \left[\bar{X} - \frac{1.96}{\sqrt{n}}, \bar{X} + \frac{1.96}{\sqrt{n}} \right] = \left[\bar{X} - 1.96\sqrt{\text{Var}(\bar{x})}, \bar{X} + 1.96\sqrt{\text{Var}(\bar{x})} \right]. \frac{\bar{X} - \mu}{\sqrt{\text{Var}(\bar{x})}} \sim \mathcal{N}(0, 1).$$

(II). Recall: The "optimal" two-sided 5% of the $\mu = \mu_0$ rejects iff $|\bar{X} - \mu_0| > \frac{1.96}{\sqrt{n}}$

$$\begin{aligned} &\Leftrightarrow \bar{X} - \mu_0 > \frac{1.96}{\sqrt{n}} \text{ or } \bar{X} - \mu_0 < -\frac{1.96}{\sqrt{n}} \\ &\Leftrightarrow \mu_0 < \bar{X} - \frac{1.96}{\sqrt{n}} \text{ or } \mu_0 > \bar{X} + \frac{1.96}{\sqrt{n}} \end{aligned}$$

Hence, the test "accepts" H_0 iff

$$\bar{X} - \frac{1.96}{\sqrt{n}} \leq \mu_0 \leq \bar{X} + \frac{1.96}{\sqrt{n}}$$

Chapter 3 Decision Rule Based Statistical Inference

3.1 Decision Rule

Given an observation $x \in X$, we want to estimate an unknown state $\theta \in S$ (not necessarily random). The θ can form x with $P_\theta(x)$. We use decision rule $\delta(x)$ to form an action (estimation of θ) $a = \hat{\theta}$.

Example:

- (1) Binary hypothesis testing (detection) when $S = \{0, 1\}$ e.g. $P_0 \sim \mathcal{N}(0, \sigma^2), P_1 \sim \mathcal{N}(\mu, \sigma^2)$
- (2) Multiple hypothesis testing (classification) when $S = \{1, 2, \dots, n\}$
- (3) (Estimation) when $S = \mathbb{R}$ e.g. $P_\theta \in \mathcal{N}(\theta, \sigma^2)$

Example 3.1 (Binary HT) For the example Binary HT, $P_0 \sim \mathcal{N}(0, \sigma^2), P_1 \sim \mathcal{N}(\mu, \sigma^2)$: decision rule $\delta : \mathbb{R} \rightarrow \{0, 1\}$

We can find a τ such that $\delta(x) = \begin{cases} 1, & x \geq \tau \\ 0, & \text{else} \end{cases} = \mathbf{1}_{x \geq \tau}$. How to choose τ ?

Type-I error probability: probability that θ is 0 but receive $\delta(x) = 1$.

$$P_I = P_0\{\delta(x) = 1\} = P_0\{x \geq \tau\} = Q\left(\frac{\tau}{\sigma}\right)$$

Type-II error probability: probability that θ is 1 but receive $\delta(x) = 0$.

$$P_{II} = P_1\{\delta(x) = 0\} = P_1(x < \tau) = Q\left(\frac{\mu - \tau}{\sigma}\right)$$

Both P_I and P_{II} depends on τ . $Q(t) = \int_t^\infty \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}} dx$

For $\tau = \frac{\mu}{2}$, $P_I = P_{II} = Q\left(\frac{\mu}{2\sigma}\right)$

Example 3.2 (Multiple HT) Consider three state $S = \{1, 2, 3\}$. We can find a τ such that $\delta(x) = \begin{cases} 1, & x < \tau_1 \\ 2, & \tau_1 \leq x \leq \tau_2 \\ 3, & x > \tau_2 \end{cases} = \mathbf{1}_{x \geq \tau}$.

Conditional Error Probabilities: probability that θ is i but receive $\delta(x) = j$ (6 types in this example)

$$P_i\{\delta(x) = j\}, \forall i \neq j$$

3.2 Maximum-Likelihood Principle (state is norandom)

Maximum-Likelihood Principle

$$\hat{\theta} = \operatorname{argmax}_{\theta \in S} P_\theta(x) = \operatorname{argmax}_{\theta \in S} \ln P_\theta(x)$$

Applied to the binary example: $P_0 \sim \mathcal{N}(0, \sigma^2)$, $P_1 \sim \mathcal{N}(\mu, \sigma^2)$.

$$P_0(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}}, P_1(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}. \ln P_0(x) = c - \frac{x^2}{2\sigma^2}, \ln P_1(x) = c - \frac{(x-\mu)^2}{2\sigma^2}.$$

Then, the rule can become

$$\hat{\theta} = \begin{cases} 0, & x^2 < (x-\mu)^2 \\ 1, & \text{else} \end{cases} = \mathbf{1}_{x^2 \geq (x-\mu)^2} = \mathbf{1}_{x \geq \frac{\mu}{2}}$$

Vector Observations

Observations $X = (x_1, x_2, \dots, x_n)$, where i.i.d. $x_i \sim P_\theta$. Then

$$P_\theta(X) = \prod_{i=1}^n P_\theta(x_i), \ln P_\theta(X) = \sum_{i=1}^n \ln P_\theta(x_i)$$

$$\ln P_0(x) = cn - \frac{\sum_{i=1}^n x_i^2}{2\sigma^2}, \ln P_1(x) = cn - \frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}.$$

Then, the rule can become

$$\hat{\theta} = \begin{cases} 0, & \sum_{i=1}^n x_i^2 < \sum_{i=1}^n (x_i - \mu)^2 \\ 1, & \text{else} \end{cases} = \mathbf{1}_{\sum_{i=1}^n x_i^2 \geq \sum_{i=1}^n (x_i - \mu)^2} = \mathbf{1}_{\bar{x} \geq \frac{\mu}{2}}$$

where $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$. Under both H_0 and H_1 , $\bar{x} \sim \mathcal{N}(0, \frac{\sigma^2}{n})$.

Then, type I error prob and type II error prob are the same

$$P_I = P_0\{\bar{x} \geq \frac{\mu}{2}\} = P_{II} = P_1\{\bar{x} < \frac{\mu}{2}\} = Q\left(\frac{\mu\sqrt{n}}{2\sigma}\right)$$

Estimation $S = \mathbb{R}$

To estimate θ when $S = \mathbb{R}$

$$\begin{aligned} & \max_{\theta \in \mathbb{R}} \sum_{i=1}^n \ln P_\theta(x_i) \\ & \Leftrightarrow \max_{\theta \in \mathbb{R}} \left[cn - \frac{\sum_{i=1}^n (x_i - \theta)^2}{2\sigma^2} \right] \\ & \Leftrightarrow \max_{\theta \in \mathbb{R}} \sum_{i=1}^n (x_i - \theta)^2 \Rightarrow \hat{\theta} = \bar{x} \end{aligned}$$

Then, with $\bar{x} \sim \mathcal{N}(\theta, \frac{\sigma^2}{n})$, the

$$MSE_\theta = \mathbb{E}_\theta (\bar{x} - \theta)^2 = \frac{\sigma^2}{n}$$

3.3 Bayesian Decision Rule (state is random)

3.3.1 Rules

Prior probability distribution of θ is $\pi(\theta)$.

Loss/cost function with action (estimation) a is $l(a, \theta)$. e.g.

1. (binary HT) Hamming/zero-one loss $l(a = \hat{\theta}, \theta) = \mathbf{1}_{a \neq \theta}$
2. (estimation) Squared error loss $l(a = \hat{\theta}, \theta) = (a - \theta)^2$; Absolute error loss $l(a, \theta) = |a - \theta|$.

Definition 3.1 (Risk)

Risk of decision rule δ on θ :

$$R(\delta, \theta) = \mathbb{E}_{X \sim \pi(\theta)} [l(\delta(X), \theta)]$$

where X are random with prob $P(\cdot | \theta)$.

Risk of decision rule δ :

$$\begin{aligned} R(\delta) &= \mathbb{E}_{\theta \sim P_\theta} [R(\delta, \theta)] \\ &= \mathbb{E}_{\theta \sim P_\theta} \mathbb{E}_{X \sim \pi(\theta)} [l(\delta(X), \theta)] \end{aligned}$$

where (X, θ) are random with joint probability distribution

$$P(X, \theta) = P(X)\pi(\theta | X)$$



Note In machine learning, we normally use y to substitute θ .

Example 3.3 (Hamming/zero-one Loss) The risk of decision $\delta(x)$ in Hamming/zero-one loss $l(a = \hat{y}, y) = \mathbf{1}_{a \neq y}$

$$\begin{aligned} R(\delta) &= \mathbb{E}(\mathbf{1}_{\delta(x) \neq y}) = \mathbb{E}[\delta(x) \neq y] \\ &= P(y = 0)P[\delta(x) \neq 0 | y = 0] + P(y = 1)P[\delta(x) \neq 1 | y = 1] \\ &= P(y = 0)P[\delta(x) = 1 | y = 0] + P(y = 1)P[\delta(x) = 0 | y = 1] \end{aligned}$$

3.3.2 Optimization Problem in Bayes Form

We want to compute the optimal rule that minimizes the risk:

$$\delta_B = \operatorname{argmin}_{\delta} R(\delta)$$

Derive Bayes rule

$$\begin{aligned} R(\delta) &= \int_x \int_{\theta} P(x, \theta) l(\delta(x), \theta) d\theta dx \\ &= \int_x P(x) \int_{\theta} \pi(\theta | x) l(\delta(x), \theta) d\theta dx \end{aligned}$$

δ_B is given by solving the optimization problem:

$$\min_{\delta} \int_x P(x) \int_{\theta} \pi(\theta | x) l(\delta(x), \theta) d\theta dx$$

Hence,

Proposition 3.1

$\delta_B = \operatorname{argmin}_{\delta} R(\delta)$ can be transformed into optimization problems for each $x \in S$

$$\min_{\delta(x)} \int_{\theta} \pi(\theta | x) l(\delta(x), \theta) d\theta$$



The problem becomes to compute $\pi(\theta | x)$, which is computed by

$$\pi(\theta | x) = \frac{\pi(\theta)P(x | \theta)}{P(x)}$$

Example 3.4 (Square Loss) Consider a Bernoulli Distribution with parameter θ . We have data $\{Z_1, \dots, Z_n\}$ and we want to predict the next sample Z .

Note that Z is what we want to predict, so the Sqaure error loss given estimation t is

$$\begin{aligned} R(t, Z; \theta) &= \mathbb{E}_{\theta}[(t - Z)^2] \\ &= (t - \theta)^2 + \theta(1 - \theta) \end{aligned}$$

Then, the optimal estimation is $t^* = \theta$. And we say the minimial risk is **oracle risk**

$$R(t^*, Z | \theta) = \theta(1 - \theta)$$

However, we don't have θ . What if we use the MLE $t(\{Z_1, \dots, Z_n\}) = \hat{\theta}_{\text{MLE}} = \bar{Z} = \frac{1}{n} \sum_{i=1}^n nZ_i$ instead?

$$\begin{aligned} R(\hat{\theta}_{\text{MLE}}, Z | \theta) &= \mathbb{E}_{\theta}[(\bar{Z} - Z)^2] \\ &= \mathbb{E}_{\theta}[(\bar{Z} - \theta)^2] + \theta(1 - \theta) \\ &= \underbrace{\frac{\theta(1 - \theta)}{n}}_{\text{Sample Uncertainty}} + \underbrace{\theta(1 - \theta)}_{\text{Oracle Risk}} \end{aligned}$$

The avenge risk with prior belief $P(\theta) = 1_{\{\theta \in [0, 1]\}}$ ($\theta \sim \text{Beta}(1, 1) = \text{Unif}[0, 1]$) is

$$R(t) = \int_{\theta} R(t, Z | \theta) \pi(\theta | \{Z_1, \dots, Z_n\}) d\theta$$

$\pi(\theta | \{Z_1, \dots, Z_n\})$ is the posterior beliefs about the θ :

$$\pi(\theta | \{Z_1, \dots, Z_n\}) = \frac{f(\{Z_1, \dots, Z_n\} | \theta)P(\theta)}{\int f(\{Z_1, \dots, Z_n\} | \theta')P(\theta')d\theta'}$$

As $Z_i \sim \text{Bernoulli}(\theta)$, we have

$$\theta | \{Z_1, \dots, Z_n\} \sim \text{Beta} \left(\sum_{i=1}^n Z_i + 1, n - \sum_{i=1}^n Z_i + 1 \right)$$

which has mean $\mathbb{E}[\theta | \{Z_1, \dots, Z_n\}] = \frac{\sum_{i=1}^n Z_i + 1}{n+2} = \frac{\hat{\theta}_{MLE} + \frac{1}{2}}{1 + \frac{2}{n}}$. Then,

$$\begin{aligned} R^*(t) &= \min_t \int_\theta R(t, Z | \theta) \pi(\theta | \{Z_1, \dots, Z_n\}) d\theta \\ &= \min_t \int_\theta (t - \theta)^2 \pi(\theta | \{Z_1, \dots, Z_n\}) d\theta + \int_\theta \theta(1 - \theta) \pi(\theta | \{Z_1, \dots, Z_n\}) d\theta \\ &\Rightarrow t^* = \mathbb{E}[\theta | \{Z_1, \dots, Z_n\}] \end{aligned}$$

3.3.3 Maximum A Posteriori (MAP) Decision Rule (Binary example)

Example 3.5 Hamming/zero-one loss $l(a, y) = \mathbf{1}_{a \neq y}$

Maximum A Posteriori (MAP) Decision Rule:

Optimization problem is

$$\begin{aligned} \delta(x) &= \operatorname{argmin}_a \sum_{y=0,1} \pi(y|x) \mathbf{1}_{a \neq y} dy = \operatorname{argmax}_{y \in \{0,1\}} \pi(y|x) \\ &\Rightarrow \sum_{y=0,1} \pi(y|x) \mathbf{1}_{\delta(x) \neq y} dx = \min_a \sum_{y=0,1} \pi(y|x) \mathbf{1}_{a \neq y} dy = \min\{\pi(1|x), \pi(0|x)\} \end{aligned}$$

Likelihood ratio: $L(x) = \frac{P_1(x)}{P_0(x)}$

Likelihood ratio test: threshold $\tau = \frac{\pi(0)}{\pi(1)}$. If $L(x) > \tau$ accept H_1 (equivalent to $P_1(x)\pi(1) > P_0(x)\pi(0)$ which is also equivalent to comparing $\pi(y|x)$).

In this rule the whole optimization problem also goes to

$$\begin{aligned} R(\delta_{MAP}) &= \int_x P(x) \sum_{y=0,1} \pi(y|x) \mathbf{1}_{\delta(x) \neq y} dx \\ &= \int_x P(x) \min\{\pi(1|x), \pi(0|x)\} dx \end{aligned}$$

3.3.4 Minimum Mean Squared Error (MMSE) Rule (\mathbb{R}^n example)

Example 3.6 (Estimation) Squared error loss $l(a, y) = (a - y)^2$.

Minimum Mean Squared Error (MMSE) Rule:

Optimization problem is $\delta(x) = \operatorname{argmin}_a \int_y \pi(y|x) (a - y)^2 dy$

$$\begin{aligned} 0 &= \int_y \pi(y|x) (\delta_B(x) - y) dy = \delta_B(x) - \mathbb{E}[Y|X=x] \\ &\Rightarrow \delta_B(x) = \mathbb{E}[Y|X=x] \end{aligned}$$

which is called **conditional mean estimation**.

In this rule the whole optimization problem also goes to

$$R(\delta_{MMSE}) = \int_x P(x) \int_y \pi(y|x)(y - \mathbb{E}[Y|X=x])^2 dy dx = \mathbb{E}_X Var[Y|X=x]$$

Gaussian case: If $X \in \mathbb{R}^n$ and (Y, X) are jointly Gaussian, then the conditional mean is a linear function of x , also called linear MMSE estimator.

$$\mathbb{E}[Y|X=x] = \mathbb{E}[Y] + Cov(Y, X)Cov(X)^{-1}(x - \mathbb{E}[X])$$

and the posterior risk is independent of x :

$$Var[Y|X=x] = Var[Y] - Cov(Y, X)Cov(X)^{-1}Cov(X, Y)$$

Note: MMSE estimator coincides with the MAP estimator for Gaussian Variables.

3.4 Comparison

Maximum-Likelihood Principle (state is nonrandom):

$$\delta_{ML}(x) = \operatorname{argmax}_y P_y(x)$$

Maximum A Posteriori (MAP) Decision Rule (state is random):

$$\delta_{MAP}(x) = \operatorname{argmax}_y \pi(y|x) = \operatorname{argmax}_y \{\pi(y|x), P_y(x)\}$$

Chapter 4 Bootstrap

Bootstrap is a procedure to compute properties of an estimator by random re-sampling with replacement from the data. It was first introduced by Efron (1979).

Suppose we have i.i.d. sample $\vec{Y} = (Y_1, Y_2, \dots, Y_n)$ taken i.i.d. from a distribution with cdf F and we want to compute a statistic θ of the distribution using an estimator $\hat{\theta}_n(\vec{Y})$. The distribution of the statistic θ has cdf G . While the estimator $\hat{\theta}_n(\vec{Y})$ may not be optimal in any sense, it is often the case that $\hat{\theta}_n(\vec{Y})$ is consistent in probability, i.e., $\hat{\theta}_n(\vec{Y}) \xrightarrow{P} \theta$ as $n \rightarrow \infty$. We want to analyze the performance of the estimotor $\hat{\theta}_n(\vec{Y})$ in terms of the following quantities:

(1). Bias:

$$\text{Bias}(\hat{\theta}_n) = \mathbb{E}_{\theta}[\hat{\theta}_n(\vec{Y})] - \theta$$

(2). Variance:

$$\text{Var}(\hat{\theta}_n) = \mathbb{E}_{\theta}[\hat{\theta}_n^2(\vec{Y})] - \mathbb{E}_{\theta}^2[\hat{\theta}_n(\vec{Y})]$$

(3). CDF:

$$G_n(t) = P(\hat{\theta}_n(\vec{Y}) < t), \forall t$$

4.1 Traditional Monte-Carlo Approach

Generate k vectors $\vec{Y}^{(i)}$, $i = 1, 2, \dots, k$ (total kn random variables)

(1). Bias:

$$\widehat{\text{Bias}}(\hat{\theta}_n) = \frac{1}{k} \sum_{j=1}^k \hat{\theta}_n(\vec{Y}^{(j)}) - \theta$$

By the strong law of large number, the mean $\frac{1}{k} \sum_{j=1}^k \hat{\theta}_n(\vec{Y}^{(j)})$ converges almost surely to the expected value $\mathbb{E}_{\theta}[\hat{\theta}_n(\vec{Y})]$, so $\widehat{\text{Bias}}(\hat{\theta}_n) \xrightarrow{a.s.} \text{Bias}(\hat{\theta}_n)$.

(2). Variance:

$$\widehat{\text{Var}}(\hat{\theta}_n) = \frac{1}{k} \sum_{j=1}^k \hat{\theta}_n^2(\vec{Y}^{(j)}) - \left(\frac{1}{k} \sum_{j=1}^k \hat{\theta}_n(\vec{Y}^{(j)}) \right)^2$$

Still by the strong law of large number, the mean $\frac{1}{k} \sum_{j=1}^k \hat{\theta}_n(\vec{Y}^{(j)})$ converges almost surely to the expected value $\mathbb{E}_{\theta}[\hat{\theta}_n(\vec{Y})]$ and the mean $\frac{1}{k} \sum_{j=1}^k \hat{\theta}_n^2(\vec{Y}^{(j)})$ converges almost surely to the expected value $\mathbb{E}_{\theta}[\hat{\theta}_n^2(\vec{Y})]$, so $\widehat{\text{Var}}(\hat{\theta}_n) \xrightarrow{a.s.} \text{Var}(\hat{\theta}_n)$.

(3). Empirical Distribution Function (CDF):

$$\hat{G}_n(t) = \frac{1}{k} \sum_{j=1}^k \mathbf{1}\{\hat{\theta}_n(\vec{Y}^{(j)}) < t\}, \forall t$$

By law of large numbers, we have $\hat{G}_n(x) \xrightarrow{a.s.} G_n(x), \forall t \in \mathbb{R}$ as $k \rightarrow \infty$.

By Glivenko-Cantelli Theorem, we have $\sup_{t \in \mathbb{R}} |\hat{G}_n(x) - G_n(x)| \xrightarrow{a.s.} 0$ as $k \rightarrow \infty$. (Stronger result).

4.2 Bootstrap (When data is not enough)

Suppose we only have data $\vec{Y} = (Y_1, \dots, Y_n)$ and we can't draw new samples from the real distribution anymore. We reuse Y_1, \dots, Y_n to obtain resamples $\vec{Y}^* = (Y_1^*, \dots, Y_n^*)$ (drawing from $\{Y_1, \dots, Y_n\}$ uniformly, equivalently drawing from the empirical distribution with cdf $F_n(y) \triangleq \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{Y_i = y\}$). We get k resamples, denoted by $\vec{Y}^{*(1)}, \dots, \vec{Y}^{*(k)}$.

1. Bias:

$$\text{Bias}^*(\hat{\theta}_n) \triangleq \frac{1}{k} \sum_{j=1}^k \hat{\theta}_n(\vec{Y}^{*(j)}) - \theta$$

2. Variance:

$$\text{Var}^*(\hat{\theta}_n) \triangleq \frac{1}{k} \sum_{j=1}^k \hat{\theta}_n^2(\vec{Y}^{*(j)}) - \left(\frac{1}{k} \sum_{j=1}^k \hat{\theta}_n(\vec{Y}^{*(j)}) \right)^2$$

3. CDF:

$$\hat{G}_n^*(t) = \frac{1}{k} \sum_{j=1}^k \mathbf{1}_{\hat{\theta}_n(\vec{Y}^{*(j)}) < t}, \forall t$$



Note $\hat{G}_n^*(t)$ may not always converge to G_n as $n \rightarrow \infty$.

Example 4.1 (Bootstrap Fail Example) Suppose $Y \sim$ i.i.d. $[0, \theta]$ and consider the estimator $\hat{\theta}_n(\vec{Y}) = \max_i Y_i \triangleq Y_{(n)}$. Then, for all $t \geq 0$,

$$G_n(t) \rightarrow 1 - e^{-\frac{t}{\theta_F}} \text{ as } n \rightarrow \infty$$

But for all $t \geq 0$,

$$\hat{G}_n^*(t) \geq P_{F_n}(Y_{(n)} = Y_{(n)}^*) = 1 - (1 - \frac{1}{n})^n \rightarrow 1 - e^{-1} \text{ as } n \rightarrow \infty$$

4.3 Residual Bootstrap (for problem with not i.i.d. data)

The bootstrap principle is quite general and may also be used in problems where the data $Y_i, 1 \leq i \leq n$, **are not i.i.d.**

4.3.1 Example: Linear

Consider the model

$$Y_i = a + bs_i + Z_i, \quad i = 1, 2, \dots, n$$

where $\theta = (a, b)$ is the parameter to be estimated, $\vec{s} = (s_1, \dots, s_n)$ is a known signal, and $Z_i \sim \mathcal{N}(0, \sigma^2)$ (i.i.d.).

The Linear Least Square Estimator is

$$(\hat{a}_n, \hat{b}_n) = \operatorname{argmin}_{(a,b)} \sum_{i=1}^n (Y_i - a - bs_i)^2$$

Given \vec{Y} and estimator $\hat{\theta}_n = (\hat{a}_n, \hat{b}_n)$, define the residual errors (not i.i.d.)

$$E_i = Y_i - \hat{a}_n - \hat{b}_n s_i \approx Z_i$$

Then, we use bootstrap to generate k resamples of $\vec{E} = (E_1, E_2, \dots, E_n)$.

For $j = 1, \dots, k$, do the following:

1. Obtain $\vec{E}^{*(j)}$ by uniformly resampling from \vec{E} .
2. Compute pseudo-data $Y_i^{*(j)} = \hat{a}_n + \hat{b}_n s_i + E_i^{*(j)}$ for $1 \leq i \leq n$.
3. Compute LS estimator to the pseudo-data

$$\hat{\theta}_n^{(j)} = (\hat{a}_n^{(j)}, \hat{b}_n^{(j)}) = \operatorname{argmin}_{(a,b)} \sum_{i=1}^n (Y_i^{*(j)} - a - bs_i)^2$$

Then, we can evaluate bias

$$\widehat{Bias} = \frac{1}{k} \sum_{j=1}^k \hat{\theta}_n^{(j)} - \theta$$

4.3.2 Example: Nonlinear Markov Process

Consider the model $Y_i = F_\theta(Y_{i-1}) + Z_i$, where $Z_i \sim \mathcal{N}(0, \sigma^2)$ (i.i.d.) for $i = 1, 2, \dots, n$

Parameter $\theta = (a, b)$. Linear Least Square Estimator:

$$\hat{\theta}_n(\vec{Y}) = \operatorname{argmin}_{\theta} \sum_{i=1}^n (Y_i - F_\theta(Y_{i-1}))^2$$

Given \vec{Y} , the residual (not i.i.d.)

$$E_i = Y_i - \hat{a}_n - F_{\hat{\theta}_n}(Y_{i-1}) \approx Z_i$$

Generate k resamples of $\vec{E} = (E_1, E_2, \dots, E_n)$

\Rightarrow obtain $\vec{E}^{*(1)}, \vec{E}^{*(2)}, \dots, \vec{E}^{*(k)}$ by resampling

\Rightarrow Fix $Y_0^{*(j)} = Y_0$, compute pseudo-data $Y_i^{*(j)} = F_{\hat{\theta}_n}(Y_{i-1}^{*(j)}) + E_i^{*(j)}$

\Rightarrow Compute LS estimator

$$\hat{\theta}_n^{(j)} = \operatorname{argmin}_{(a,b)} \sum_{i=1}^n (Y_i^{*(j)} - F_{\hat{\theta}_n}(Y_{i-1}^{*(j)}))^2$$

⇒ Evaluate bias

$$\widehat{Bias} = \frac{1}{k} \sum_{j=1}^k \hat{\theta}_n^{(j)} - \theta$$

4.4 Posterior Simulation / Bayesian (Weighted) Bootstrap

Assumption *Bootstrap makes a strong assumption: The data is discrete and values not seen in the data are impossible.*

Consider $Z \in \mathbb{Z} = \{z_1, \dots, z_J\}$ with parameter $\vec{\theta} = \{\theta_1, \dots, \theta_J\} \in \Theta = \mathbb{S}^{J-1} = \{\vec{\theta} \in \mathbb{R}^J : \sum_{j=1}^J \theta_j = 1, \theta_j \geq 0, j = 1, \dots, J\}$ such that $P(Z = z_j \mid \vec{\theta}) = \theta_j$.

Given a sample $\vec{Z} = (Z_1, \dots, Z_N)$. Define $N_j = \sum_{i=1}^N \mathbf{1}\{Z_i = z_j\}, j = 1, 2, \dots, J$, the number of observations that have value z_j . Then, the conditional pmf of $\vec{Z} \mid \vec{\theta}$ is

$$f(\vec{Z} \mid \vec{\theta}) = \prod_{j=1}^J \theta_j^{N_j}$$

Definition 4.1 (Steps to estimate β by Bayesian Bootstrap)

- (1). We have prior $\pi(\vec{\theta})$.
- (2). Given \vec{Z} , calculate posterior distribution $\pi(\vec{\theta} \mid \vec{Z})$.
- (3). Draw samples $\vec{\theta}^{(t)}, t = 1, \dots, T$ from $\pi(\vec{\theta} \mid \vec{Z})$.
- (4). Then compute $\frac{1}{T} \sum_{t=1}^T \hat{\beta}(\vec{\theta}^{(t)})$.



4.4.1 Dirichlet Distribution Prior

A convenient way to assign the prior distribution of $\vec{\theta}$ over Θ is to use Dirichlet distribution.

Definition 4.2 (Dirichlet Distribution)

A **Dirichlet distribution** with parameters $\vec{\alpha} = (\alpha_1, \dots, \alpha_J), J \geq 2$. It allocates mass on $\vec{\theta}$ over Θ ,

$$\pi(\vec{\theta}) = \frac{\Gamma(\sum_{j=1}^J \alpha_j)}{\sum_{j=1}^N \Gamma(\alpha_j)} \prod_{j=1}^J \theta_j^{\alpha_j - 1}$$

where $\Gamma(z) \triangleq \int_0^\infty t^{z-1} e^{-t} dt$ is Gamma function (if z is positive integer, $\Gamma(z) = (z-1)!$).



Note Dirichlet distribution generalizes Beta distribution.



Now let's use Dirichlet distribution with parameters $\vec{\alpha} = (\alpha_1, \dots, \alpha_J)$ to estimate $\mathbb{E}[\vec{\theta} \mid \vec{Z}]$.

As $f(\vec{Z} \mid \vec{\theta}) = \prod_{j=1}^J \theta_j^{N_j}$, we can compute the posterior beliefs

$$\pi(\vec{\theta} \mid \vec{Z}) = \frac{f(\vec{Z} \mid \vec{\theta}) P(\vec{\theta})}{\int f(\vec{Z} \mid \vec{\theta}') P(\vec{\theta}') d\vec{\theta}'} = \frac{\Gamma(\sum_{j=1}^J (N_j + \alpha_j))}{\sum_{j=1}^N \Gamma(N_j + \alpha_j)} \prod_{j=1}^J \theta_j^{N_j + \alpha_j - 1}$$

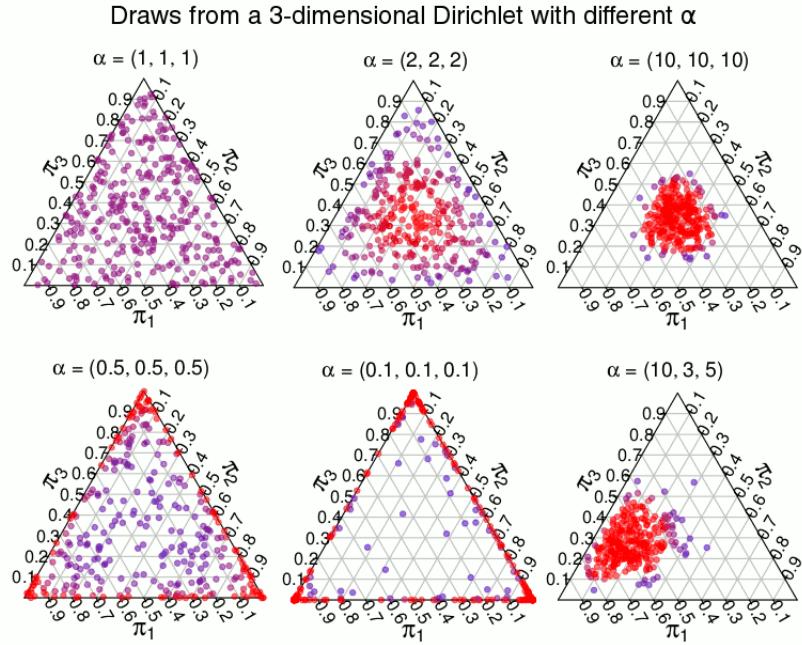


Figure 4.1: Dirichlet Distribution Examples

That is

$$\theta \mid \vec{Z} \sim \text{Dirichlet}(\vec{\alpha}), \text{ where } \bar{\alpha}_j = \alpha_j + N_j, \forall j$$

Simulate samples from Dirichlet distribution

Definition 4.3 (Simulate samples from $\text{Dirichlet}(\vec{\alpha})$)

1. Consider a series of independent Gamma random variable $w_i \sim \text{Gamma}(\alpha_i, 1), i = 1, \dots, J$;
2. Define $v_i = \frac{w_i}{\sum_{j=1}^J w_j}$;
3. We have $(v_1, \dots, v_J) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_J)$.



4.4.2 Haldane Prior

We may also begin with an uninformative prior, an improper prior, $\text{Dirichlet}(\vec{\alpha})$, where $\vec{\alpha} \rightarrow 0$. $\pi(\theta) \propto \frac{1}{\theta_1 \theta_2 \dots \theta_J}$.

Under this prior, the posterior is $\text{Dirichlet}(N_1, \dots, N_J)$, where $N_j = \sum_{i=1}^N \mathbf{1}\{Z_i = z_j\}$.

4.4.3 Linear Model Case

Each sample is $Z_i = (1, X_{1,i}, X_{2,i}, X_{3,i}, X_{4,i})$. The linear regression coefficient is $\beta = \mathbb{E}[XX']^{-1}\mathbb{E}[XY]$, and $\mathbb{E}^*[Y \mid X = x] = x'\beta$.

4.4.4 Bernoulli Case

Consider the problem of Example 3.4. Given N random sample $\{Z_1, \dots, Z_N\}$ from a Bernoulli distribution with parameter θ and the sum $\sum_{i=1}^N Z_i = S$.

Consider a series of Gamma random variable $w_i^{(t)} \sim \text{Gamma}(1, 1)$ from time $t = 1, \dots, T$. Then, we have

$$\begin{aligned}\sum_{i=1}^N w_i^{(t)} \mathbf{1}_{\{Z_i=1\}} &\sim \text{Gamma}(S, 1) \\ \sum_{i=1}^N w_i^{(t)} \mathbf{1}_{\{Z_i=0\}} &\sim \text{Gamma}(N - S, 1)\end{aligned}$$

Define $v_i^{(t)} = \frac{w_i^{(t)}}{\sum_{j=1}^N w_j^{(t)}}$. Based on the property of Gamma distribution, we have $\mathbb{E}[w_i^{(t)}] = \text{Var}[w_i^{(t)}] = 1$ and $\mathbb{E}[v_i^{(t)}] = \frac{1}{N}$.

As the relation between Gamma distribution and Beta distribution, we have

$$\frac{\text{Gamma}(S, 1)}{\text{Gamma}(S, 1) + \text{Gamma}(N - S, 1)} \sim \text{Beta}(S, N - S)$$

Hence, we can define

$$\begin{aligned}\hat{\theta}^{(t)} &= \sum_{i=1}^N v_i^{(t)} Z_i \\ &= \sum_{i=1}^N \frac{w_i^{(t)} Z_i}{\sum_{j=1}^N w_j^{(t)}} \sim \text{Beta}(S, N - S)\end{aligned}$$

which is close to the posterior beliefs in Example 3.4 and can be seen as the posterior beliefs drawn from an improper prior: $\theta \sim \text{Beta}(\epsilon, \epsilon)$, $\epsilon \rightarrow 0$, which has p.d.f. $\pi(\theta) = \frac{1}{\theta(1-\theta)}$.

We use

$$\frac{1}{T} \sum_{t=1}^T \hat{\theta}^{(t)} \approx \mathbb{E}[\theta^{(t)} | \{Z_1, \dots, Z_n\}]$$

to estimate $\mathbb{E}[\theta^{(t)} | \{Z_1, \dots, Z_n\}]$.

Chapter 5 Nonparameteric Prediction Problem

Problem

There are J non-stochastic treats $X \in \mathbb{X} \subseteq \mathbb{R}^J$, and we want to predict a related outcome $Y \in \mathbb{Y} \subseteq \mathbb{R}$.

Given a sample X_i ,

$$Y_i = m(X_i) + \sigma u_i$$

where $m(\cdot)$ is an unknown function and $u_i | X_i \sim \mathcal{N}(0, 1)$.

Goal:

- o Predict Y given a new X ;
- o Learn $m(\cdot)$.

Decision Rule

Given $\vec{X} = \{X_1, \dots, X_N\}$, we want to derive a decision rule $d(\vec{Y})$ given corresponding \vec{Y} of \vec{X} .

Define $\mathbf{m} \triangleq [m(X_1), \dots, m(X_N)]'$, its estimation is denoted by $\hat{\mathbf{m}}$, which is based on the decision rule.

Sum of Squared Residual

Proposition 5.1 (Sum of Squared Residual)

Sum of Squared Residual (SSR) of an estimation $\hat{\mathbf{m}}$ is given by

$$\mathbb{E} [\|\vec{Y} - \hat{\mathbf{m}}\|^2] = N\sigma^2 + \sum_{i=1}^N (\hat{m}(X_i) - m(X_i))^2 - 2\sigma^2 df(\hat{\mathbf{m}}) \quad (\text{SSR})$$

where the norm is $\|\vec{X}\| = \left[\sum_{i=1}^N X_i^2 \right]^{1/2}$, and the degree of freedom $df(\hat{\mathbf{m}}) = \frac{\sum_{i=1}^N \text{Cov}(Y_i, \hat{m}_i)}{\sigma^2}$.



Proof 5.1

$$\begin{aligned} \mathbb{E} [\|\vec{Y} - \hat{\mathbf{m}}\|^2] &= \mathbb{E} [\|(\vec{Y} - \mathbf{m}) + (\mathbf{m} - \hat{\mathbf{m}})\|^2] \\ &= \mathbb{E} [\|\vec{Y} - \mathbf{m}\|^2] + \mathbb{E} [\|\hat{\mathbf{m}} - \mathbf{m}\|^2] - 2\mathbb{E} [(\vec{Y} - \mathbf{m})'(\hat{\mathbf{m}} - \mathbf{m})] \\ &= N\sigma^2 + \sum_{i=1}^N (\hat{m}(X_i) - m(X_i))^2 - 2\sigma^2 df(\hat{\mathbf{m}}) \end{aligned}$$

The second equality is because

$$\begin{aligned}
 \mathbb{E} [(\vec{Y} - \mathbf{m})'(\hat{\mathbf{m}} - \mathbf{m})] &= \sum_{i=1}^N \mathbb{E} [(Y_i - m_i)(\hat{m}_i - m_i)] \\
 &= \sum_{i=1}^N \mathbb{E} [(Y_i - m_i)\hat{m}_i] \\
 &= \sum_{i=1}^N \mathbb{E} [(Y_i - m_i)(\hat{m}_i - \mathbb{E} m_i)] \\
 &= \sum_{i=1}^N \text{Cov}(Y_i, \hat{m}_i)
 \end{aligned}$$

We can represent the risk of estimation $\hat{\mathbf{m}}$ by rewriting **SSR**

$$\mathbb{E} [\|\hat{\mathbf{m}} - \mathbf{m}\|^2] = \mathbb{E} [\|\vec{Y} - \hat{\mathbf{m}}\|^2] - N\sigma^2 + 2\sigma^2 df(\hat{\mathbf{m}})$$

5.1 K-normal Means Probelm

5.1.1 Assumptions

Assumption

1. *Linear Combination:* Suppose the $m(\cdot)$ can be written as a linear combination of basis functions:

$$\mathbf{m}(X) = \sum_{k=1}^K \alpha_k g_k(X)$$

2. *Gram-Schmidt Orthonormalization* (How to process raw data):

Definition 5.1 (Gram-Schmidt Orthonormalization)

$\phi_k(X), k = 1, \dots, K$ such that

- (1). $\frac{1}{N} \sum_{i=1}^N \phi_k^2(X_i) = 1$ and
- (2). $\frac{1}{N} \sum_{i=1}^N \phi_k(X_i) \phi_k(X_j) = 0$.



Suppose the $\mathbf{m}(\cdot)$ is a linear combination of Gram-Schmidt orthonormalizations:

$$\mathbf{m}(X) = \sum_{k=1}^K \theta_k \phi_k(X) \triangleq W\boldsymbol{\theta}$$

where $W = (w(X_1), \dots, w(X_N))^T \in \mathbb{R}^{N \times K}$, $w(X_i) = (\phi_1(X_i), \dots, \phi_K(X_i))^T \in \mathbb{R}^{K \times 1}$, and $\boldsymbol{\theta} = (\theta_1, \dots, \theta_K)^T$.

That is, given X , let

$$\underbrace{\mathbf{Y}}_{N \times 1} = \underbrace{W\boldsymbol{\theta}}_{(N \times K)(K \times 1)} + \underbrace{\sigma^2 U}_{N \times 1}, \text{ where } U \sim \mathcal{N}(0, I_N)$$

5.1.2 Maximum Likelihood Estimator

Based on these assumptions, the conditional distribution is

$$\mathbf{Y} | X \sim \mathcal{N}(W\boldsymbol{\theta}, \sigma^2 I_N)$$

and the log-likelihood function is

$$l(Y | X, \boldsymbol{\theta}) = -\frac{N}{2} \ln 2\pi - N \ln \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^N (Y_i - w(X_i)^T \boldsymbol{\theta})^2$$

Then, we get the maximum likelihood estimator (MLE),

$$\hat{\boldsymbol{\theta}}_{MLE} = \left[\frac{1}{N} \sum_{i=1}^N w(X_i)w(X_i)^T \right]^{-1} \left[\frac{1}{N} \sum_{i=1}^N w(X_i)Y_i \right]$$

By the orthonormalization assumption 5.1, $\frac{1}{N} \sum_{i=1}^N w(X_i)w(X_i)^T = I_K$. Hence,

$$\hat{\boldsymbol{\theta}}_{MLE} = \frac{1}{N} \sum_{i=1}^N w(X_i)Y_i = \frac{W^T \mathbf{Y}}{N} \in \mathbb{R}^{K \times 1}$$

We can observe that it is a **conditionally unbiased** estimate of $\boldsymbol{\theta}$:

$$\mathbb{E}[\hat{\boldsymbol{\theta}}_{MLE} | X] = \frac{1}{N} \sum_{i=1}^N w(X_i) \mathbb{E}[Y_i | X] = \left(\frac{1}{N} \sum_{i=1}^N w_i w_i^T \right) \boldsymbol{\theta} = \boldsymbol{\theta}$$

The **variance** of k^{th} item of $\hat{\boldsymbol{\theta}}_{MLE}$, $Z_k \triangleq \frac{1}{N} \sum_{i=1}^N \phi_k(X_i)Y_i | X$, is

$$\text{Var}(Z_k) = \text{Var}\left(\frac{1}{N} \sum_{i=1}^N \phi_k(X_i)Y_i | X\right) = \frac{1}{N^2} \sum_{i=1}^N \phi_k^2(X_i) \text{Var}(Y_i | X) = \frac{\sigma^2}{N}$$

Hence,

$$\hat{\boldsymbol{\theta}}_{MLE} | X \sim \mathcal{N}\left(\boldsymbol{\theta}, \frac{\sigma^2}{N} I_K\right)$$

5.1.3 Risk of MLE

All in all, we estimate the $\mathbf{m}(\cdot)$ by the unbiased estimator

$$\hat{\mathbf{m}} = W\hat{\boldsymbol{\theta}}_{MLE}$$

and then the loss is given by

$$\|\hat{\mathbf{m}} - \mathbf{m}\|^2 = \|W(\hat{\boldsymbol{\theta}}_{MLE} - \boldsymbol{\theta})\|^2 = \sum_{k=1}^K (Z_k - \theta_k)^2 \triangleq L(\hat{\boldsymbol{\theta}}_{MLE}, \boldsymbol{\theta})$$

where we consider the square loss. Hence, the risk of MLE estimation is

$$\begin{aligned} R(d_{MLE}, \boldsymbol{\theta}) &= \mathbb{E}[\|\hat{\mathbf{m}} - \mathbf{m}\|^2] \\ &= \sum_{k=1}^K \mathbb{E}[(Z_k - \theta_k)^2] = \frac{K}{N} \sigma^2 \end{aligned}$$

5.1.4 James-Stein Type Estimator

MLE is a member of the class of estimators $\mathcal{L} = \{C\mathbb{Z} : C = \{c_1, \dots, c_K\}, c_k \in [0, 1]\}$. Here, we consider a estimartor in the class:

$$\begin{aligned}\mathbb{E} \left[\sum_{k=1}^K (c_k Z_k - \theta_k)^2 \right] &= \mathbb{E} \left[\sum_{k=1}^K (c_k(Z_k - \theta_k) - (1 - c_k)\theta_k)^2 \right] \\ &= \frac{\sigma^2}{N} \sum_{i=1}^K c_k^2 + \sum_{i=1}^K (1 - c_k)^2 \theta_k^2\end{aligned}$$

By F.O.C., the optimal estimator minimizing the risk is

$$c_k^* = \frac{\theta_k^2}{\frac{\sigma^2}{N} + \theta_k^2}, \quad k = 1, \dots, K$$

Then, the oracle bound is

$$\inf_{d_L} R(d_L, \boldsymbol{\theta}) = \underbrace{\frac{\sigma^2}{N} \left(\sum_{k=1}^K \frac{\theta_k^2}{\frac{\sigma^2}{N} + \theta_k^2} \right)}_{\text{oracle bound}} < \frac{K}{N} \sigma^2$$

(we can't achieve it as we don't know $\boldsymbol{\theta}$).

Is there a feasible estimator which uniformly improves upon MLE? Yes!

5.1.4.1 Stein's Unbiased Risk Estimate (SURE)

Consider sample $\mathcal{Z} \sim \mathcal{N}(\theta, \sigma^2 I_K)$, and a estimator based on \mathcal{Z} , $\hat{\theta} = \hat{\theta}(\mathcal{Z})$. Let $g(\hat{\theta}) = \hat{\theta} - \mathcal{Z}$. Then, the estimate of risk is given by

$$\hat{R}_{\text{SURE}}(\mathcal{Z}) = K\sigma^2 + 2\sigma^2 \sum_{k=1}^K \frac{\partial g_k(\mathcal{Z})}{\partial \mathcal{Z}_k} + \underbrace{\sum_{k=1}^K (\hat{\theta}_k - \mathcal{Z}_k)^2}_{\|g(\hat{\theta})\|_{\text{Frobenius}}^2}$$

It is an unbiased estimate of the mean-squared error:

$$\mathbb{E}[\hat{R}_{\text{SURE}}] = \mathbb{E}[\|\hat{\theta} - \theta\|^2]$$

Proof 5.2 (for the unbiased property)

According to (SSR),

$$\begin{aligned}\mathbb{E}[\|\hat{\theta} - \mathcal{Z}\|^2] &= \mathbb{E}[\|\hat{\theta} - \theta\|^2] + K\sigma^2 - 2 \sum_{k=1}^K \text{Cov}(\mathcal{Z}_k, \hat{\theta}_k) \\ \Rightarrow \mathbb{E}[\|\hat{\theta} - \theta\|^2] &= -K\sigma^2 + 2 \sum_{k=1}^K \text{Cov}(\mathcal{Z}_k, \hat{\theta}_k) + \mathbb{E}[\|\hat{\theta} - \mathcal{Z}\|^2]\end{aligned}$$

where

$$\begin{aligned}\text{Cov}(\mathcal{Z}_k, \hat{\theta}_k) &= \mathbb{E}[\hat{\theta}_k(\mathcal{Z}_k - \theta_k)] \\ &= \mathbb{E}[(\hat{\theta}_k - Z_k)(\mathcal{Z}_k - \theta_k)] + \mathbb{E}[(\mathcal{Z}_k - \theta_k)^2] \\ &= \mathbb{E}[g_k(\mathcal{Z})(\mathcal{Z}_k - \theta_k)] + \sigma^2\end{aligned}$$

Note $g(\hat{\theta}) = \hat{\theta} - \mathcal{Z}$.

Claim 5.1

$$\mathbb{E}[g_k(\mathcal{Z})(\mathcal{Z}_k - \theta_k)] = \sigma^2 \mathbb{E}\left[\frac{\partial g_k(\mathcal{Z})}{\partial \mathcal{Z}_k}\right]$$



$$\begin{aligned}\sigma^2 \mathbb{E}[\nabla_Z g(\mathbf{Z})] &= \sigma^2 \int_z f_Z(z) \nabla_z g(z) dz \\ &= \sigma^2 \int_z f_Z(z) dg(z) \\ &= \sigma^2 \left(f_Z(z) g(z) \Big|_{\partial \mathbb{R}^K} - \int_z g(z) \frac{\partial f_Z(z)}{\partial z} dz \right) \\ &= \sigma^2 \int_z g(z) \left(\frac{1}{\sigma^2} (z - \theta)' f_Z(z) \right) dz \\ &= \int_z f_Z(z) g(z) (z - \theta)' dz \\ &= \mathbb{E}[g(\mathbf{Z})(\mathbf{Z} - \theta)']\end{aligned}$$

Hence,

$$\mathbb{E}[\|\hat{\theta} - \theta\|^2] = K\sigma^2 + 2\sigma^2 \sum_{k=1}^K \mathbb{E}\left[\frac{\partial g_k(\mathcal{Z})}{\partial \mathcal{Z}_k}\right] + \mathbb{E}[\|\hat{\theta} - \mathcal{Z}\|^2] = \mathbb{E}[\hat{R}_{\text{SURE}}]$$

5.1.4.2 James and Stein Estimator

Note: Here, we consider $\mathcal{Z} \sim \mathcal{N}(\theta, \frac{\sigma^2}{N} I_K)$

Theorem 5.1 (James and Stein (1961))

$$\hat{\theta}_{JS}(\mathcal{Z}) = \left(1 - \frac{(K-2)\sigma^2}{\mathcal{Z}'\mathcal{Z}} \frac{N}{N}\right) \mathcal{Z}$$



We have, $g_{JS}(\mathcal{Z}) = -\frac{(K-2)\sigma^2}{\mathcal{Z}'\mathcal{Z}} \frac{N}{N} \mathcal{Z}$ and

$$\sum_{k=1}^K \frac{\partial g_k(\mathcal{Z})}{\partial \mathcal{Z}_k} = -\frac{(K-2)\sigma^2}{\mathcal{Z}'\mathcal{Z}} \frac{N}{N} \sum_{k=1}^K \left(1 - \frac{2\mathcal{Z}_k^2}{\mathcal{Z}'\mathcal{Z}}\right) = -\frac{(K-2)^2\sigma^2}{\mathcal{Z}'\mathcal{Z}} \frac{N}{N}$$

Hence, the corresponding SURE is

$$\begin{aligned}\hat{R}_{\text{SURE}}(\mathcal{Z}) &= \frac{K}{N}\sigma^2 - \frac{2\sigma^2}{N} \frac{(K-2)^2}{\mathcal{Z}'\mathcal{Z}} \frac{\sigma^2}{N} + \frac{(K-2)^2}{(\mathcal{Z}'\mathcal{Z})^2} \left(\frac{\sigma^2}{N}\right)^2 \sum_{k=1}^K \mathcal{Z}_k^2 \\ &= \frac{K}{N}\sigma^2 - \frac{(K-2)^2}{\mathcal{Z}'\mathcal{Z}} \frac{\sigma^4}{N^2}\end{aligned}$$

Then,

$$\begin{aligned}R(\hat{\theta}_{JS}, \theta) &= \mathbb{E}[\hat{R}_{\text{SURE}}(\mathcal{Z})] \\ &= \frac{K}{N}\sigma^2 - (K-2)^2 \frac{\sigma^4}{N^2} \mathbb{E}\left[\frac{1}{\mathcal{Z}'\mathcal{Z}}\right]\end{aligned}\tag{RJS}$$

As $\mathcal{Z}_k \sim \mathcal{N}(\theta_k, \frac{\sigma^2}{N})$, $\mathcal{Z}'\mathcal{Z} = \sum_{k=1}^K \mathcal{Z}_k^2 \sim \frac{\sigma^2}{N}V$ such that $V \sim \chi_{K+2W}^2$ where $W \sim \text{Poisson}(\frac{\rho}{2})$ and $\rho = N \sum_{k=1}^K \frac{\theta_k^2}{\sigma^2}$. So,

$$\begin{aligned}\mathbb{E}\left[\frac{1}{\mathcal{Z}'\mathcal{Z}}\right] &= \frac{N}{\sigma^2} \mathbb{E}\left[\frac{1}{V}\right] \\ &= \frac{N}{\sigma^2} \mathbb{E}\left[\frac{1}{K-2+2W}\right] \text{ (by the identity of chi-square distribution)} \\ &\geq \frac{N}{\sigma^2} \frac{1}{K-2+\rho} \text{ (by Jensen's inequality)} \\ &= \frac{1}{(K-2)\frac{\sigma^2}{N} + \|\theta\|_2^2}\end{aligned}$$

Substitute it into RJS,

$$R(\hat{\theta}_{JS}, \theta) \leq \frac{K}{N}\sigma^2 - \frac{(K-2)^2 \frac{\sigma^4}{N^2}}{(K-2)\frac{\sigma^2}{N} + \|\theta\|_2^2}$$

Hence,

$$R(\hat{\theta}_{MLE}, \theta) - R(\hat{\theta}_{JS}, \theta) \geq \frac{(K-2)^2 \frac{\sigma^4}{N^2}}{(K-2)\frac{\sigma^2}{N} + \|\theta\|_2^2} \geq 0$$

which shows that $\hat{\theta}_{JS}$ works better than $\hat{\theta}_{MLE}$ (ML is inadmissible) under squared error loss.

5.1.4.3 A more general form of estimator $\mathcal{L} = \{C\mathbf{Z} : C = \text{diag}\vec{c}, \vec{c} \in [0, 1]^K\}$

Consider a new estimator $\hat{\theta} \in \mathcal{L} = \{C\mathbf{Z} : C = \text{diag}\vec{c}, \vec{c} \in [0, 1]^K\}$. The SURE is

$$\hat{R}_{\text{SURE}}(\mathbf{Z}, \vec{c}) = \frac{\sigma^2}{N} \sum_{k=1}^K c_k^2 + \sum_{k=1}^K \left(\mathcal{Z}_k^2 - \frac{\sigma^2}{N}\right) (1 - c_k)^2$$

Empirical Risk Minimization: taking F.O.C.

$$\hat{c}_k = \left(1 - \frac{\sigma^2}{N} \frac{1}{\mathcal{Z}_k^2}\right), \quad k = 1, \dots, K$$

Chapter 6 Linear Predictors / Regression

6.1 Best Linear Predictor

Consider a prediction problem that the distribution $F_{X,Y}$ is known, we observe $X = \begin{pmatrix} 1 \\ R \end{pmatrix} \in \mathbb{R}^{K \times 1}$ and predict $Y \in \mathbb{R}$. Only linear functions of X are allowed $\mathcal{L} = \{X'b : b \in \mathbb{R}^K\}$. We use square experience loss $(Y - X'b)^2$. We want to minimize Risk (mean squared error)

$$\mathbb{E}_{X,Y}[(Y - X'b)^2] = \int_{x,y} (y - x'b)^2 f_{x,y}(x, y) dx dy$$

Assumption Following inference is based on assumptions:

- (i). $\mathbb{E}[Y^2] < \infty$;
- (ii). $\mathbb{E}[\|X\|^2] < \infty$ (Frobenius norm);
- (iii). $\mathbb{E}[(\alpha' X)^2] > 0$ for any non-zero $\alpha \in \mathbb{R}^K$.

Let $\beta_0 = \arg \min_{b \in \mathbb{R}^k} \mathbb{E}_{X,Y}[(Y - X'b)^2]$. By the F.O.C.

$$\mathbb{E}[X(Y - X'\beta_0)] = 0$$

$$\mathbb{E}[XY] - \mathbb{E}[XX']\beta_0 = 0$$

$$\mathbb{E}[XY] = \underbrace{\mathbb{E}[XX']}_{non-singular} \beta_0$$

$$\beta_0 = \mathbb{E}[XX']^{-1}\mathbb{E}[XY]$$

Proposition 6.1 (Best Linear Predictor)

Hence, the mean-squared error minimizing linear predictor of Y given X is

$$\mathbb{E}^*[Y|X] = X'\beta_0, \text{ where } \beta_0 = \mathbb{E}[XX']^{-1}\mathbb{E}[XY]$$



$$\mathbb{E}_{X,Y}[X(\underbrace{Y - X'\beta_0}_u)] = \begin{pmatrix} \mathbb{E}[u] \\ \mathbb{E}[uR] \end{pmatrix} = \mathbf{0}$$

Hence, we have $\mathbb{E}[u] = 0$, then $\mathbb{E}[uR] = 0 = \text{Cov}(u, R)$.

Lemma 6.1

$\mathbb{E}[u] = \mathbb{E}[uR] = \text{Cov}(u, R) = 0$, where $u = Y - \mathbb{E}^*[Y|X]$.



If $u > 0$, it is underpredicting and if $u < 0$, it is overpredicting.

Result 1 (ure Partitioned Inverse Formula)

When we separate the constant term from other variables, we can write the Best Linear Predictor as:

Proposition 6.2 (Best Linear Predictor (ure Partitioned Inverse Formula))

$$X = \begin{pmatrix} 1 \\ R \end{pmatrix}, \beta_0 = \begin{pmatrix} \alpha_0 \\ \beta_* \end{pmatrix}, \mathbb{E}[XX']^{-1} = \begin{bmatrix} 1 & \mathbb{E}[R]' \\ \mathbb{E}[R] & \mathbb{E}[RR'] \end{bmatrix}^{-1}, \mathbb{E}[XY] = \begin{pmatrix} \mathbb{E}[Y] \\ \mathbb{E}[RY] \end{pmatrix}. \text{ Then,}$$

$$\alpha_0 = \mathbb{E}[Y] - \mathbb{E}[R]'\beta_*$$

$$\beta_* = \underbrace{\text{Var}(R)^{-1}}_{(K-1) \times (K-1)} \times \underbrace{\text{Cov}(R, Y)}_{(K-1) \times 1}$$



6.2 Convergence of OLS

6.2.1 Approximation

OLS Fit is

$$\hat{\beta} = \left[\frac{1}{N} \sum_{i=1}^N X_i X_i' \right]^{-1} \left[\frac{1}{N} \sum_{i=1}^N X_i Y_i \right]$$

Theorem 6.1 (Weak Law of Large Numbers (wLLN))

The weak law of large numbers (also called Khinchin's law) states that the sample average converges in probability towards the expected value.

$$\bar{X}_n \xrightarrow{P} \mu \quad \text{when } n \rightarrow \infty.$$

That is, for any positive number ε ,

$$\lim_{n \rightarrow \infty} \Pr(|\bar{X}_n - \mu| < \varepsilon) = 1.$$



1. By LLN: $\frac{1}{N} \sum_{i=1}^N X_i Y_i \xrightarrow{P} \mathbb{E}[XY]$
2. By LLN and $f(X) = X^{-1}$ is continuous, $\left[\frac{1}{N} \sum_{i=1}^N X_i X_i' \right] \xrightarrow{P} \mathbb{E}[XX']^{-1}$
3. Hence,

$$\hat{\beta} = \left[\frac{1}{N} \sum_{i=1}^N X_i X_i' \right]^{-1} \left[\frac{1}{N} \sum_{i=1}^N X_i Y_i \right] \xrightarrow{P} \mathbb{E}[XX']^{-1} \mathbb{E}[XY] = \beta_0$$

Theorem 6.2 (Central Limit Theorem (CLT))

$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \xrightarrow{D} N(0, 1) \text{ when } n \rightarrow \infty$$

Z converges in distribution to $N(0, 1)$ as $n \rightarrow \infty$

$$(\text{converges in distribution: } P(\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \leq a) \rightarrow \frac{1}{\sqrt{2\pi}} \int_{-\infty}^a e^{-\frac{x^2}{2}} dx)$$



Application to OLS: Let $u = Y - X'\beta_0$. Then,

$$\begin{aligned}\hat{\beta} &= \left[\frac{1}{N} \sum_{i=1}^N X_i X'_i \right]^{-1} \left[\frac{1}{N} \sum_{i=1}^N X_i Y_i \right] \\ &= \left[\frac{1}{N} \sum_{i=1}^N X_i X'_i \right]^{-1} \left[\frac{1}{N} \sum_{i=1}^N X_i (u_i + X'_i \beta_0) \right] \\ &= \beta_0 + \left[\frac{1}{N} \sum_{i=1}^N X_i X'_i \right]^{-1} \left[\frac{1}{\sqrt{N}} \sum_{i=1}^N X_i u_i \right]\end{aligned}$$

Then,

$$\sqrt{N}(\hat{\beta} - \beta_0) = \left[\frac{1}{N} \sum_{i=1}^N X_i X'_i \right]^{-1} \left[\frac{1}{\sqrt{N}} \sum_{i=1}^N X_i u_i \right]$$

1. By LLN, $\left[\frac{1}{N} \sum_{i=1}^N X_i X'_i \right]^{-1} \xrightarrow{P} \mathbb{E}[X X']^{-1} \triangleq \Gamma_0^{-1}$.

2. By CLT, $\left[\frac{1}{\sqrt{N}} \sum_{i=1}^N X_i u_i \right] \sim \mathcal{N}(0, \Omega_0)$, where

$$\Omega_0 = \text{Var}[X_i u_i] = \mathbb{E}[\|X_i u_i\|^2] = \mathbb{E}[\|x_i\|^2 u_i^2] \leq (\mathbb{E}[\|x_i\|^4])^{\frac{1}{2}} \mathbb{E}[u_i^4]^{\frac{1}{2}}$$

Hence,

$$\sqrt{N}(\hat{\beta} - \beta_0) \xrightarrow{D} N(0, \Gamma_0^{-1} \Omega_0 \Gamma_0^{-1})$$

The estimation of Γ_0 and Ω_0 :

$$\begin{aligned}\hat{\Gamma} &= \frac{1}{N} \sum_{i=1}^N X_i X'_i \\ \hat{\Omega} &= \frac{1}{N} \sum_{i=1}^N X_i \hat{u}_i \hat{u}_i' X'_i, \quad \text{where } \hat{u}_i = Y_i - X'_i \hat{\beta}\end{aligned}$$

We have

$$\hat{\Gamma}^{-1} \hat{\Omega} \hat{\Gamma}^{-1} \xrightarrow{P} \Gamma_0^{-1} \Omega_0 \Gamma_0^{-1}$$

Then,

$$\hat{\beta} \xrightarrow{\text{approx}} N\left(\beta_0, \frac{\hat{\Gamma}^{-1} \hat{\Omega} \hat{\Gamma}^{-1}}{N}\right)$$

6.2.2 Testing and Confidence Interval

Let $\hat{\Lambda} = \hat{\Gamma}^{-1}\hat{\Omega}\hat{\Gamma}^{-1}$, $\Lambda = \Gamma_0^{-1}\Omega_0\Gamma_0^{-1}$, $\sqrt{N}(\hat{\beta}_k - \beta_k) \xrightarrow{D} N(0, \Lambda_{kk})$. Hence,

$$T_N \triangleq \sqrt{N}\Lambda_{kk}^{-\frac{1}{2}} (\hat{\beta}_k - \beta_k) \xrightarrow{D} N(0, 1)$$

Consider the event $A = \mathbf{1}\{|T_N| \leq 1.96\}$. We have

$$\Pr(A = 1) = \Phi(1.96) - \Phi(-1.96) = 0.95$$

Specifically,

$$\begin{aligned} A &= \mathbf{1}\{|T_N| \leq 1.96\} \\ &= \mathbf{1}\left\{\hat{\beta}_k - 1.96\frac{\Lambda_{kk}^{\frac{1}{2}}}{\sqrt{N}} \leq \beta_k \leq \hat{\beta}_k + 1.96\frac{\Lambda_{kk}^{\frac{1}{2}}}{\sqrt{N}}\right\} \end{aligned}$$

The ‘‘Random Interval’’ is

$$\left[\hat{\beta}_k - 1.96\frac{\Lambda_{kk}^{\frac{1}{2}}}{\sqrt{N}}, \hat{\beta}_k + 1.96\frac{\Lambda_{kk}^{\frac{1}{2}}}{\sqrt{N}}\right]$$

Testing Linear Restrictions

Let $\theta = H\beta$, where H is $p \times k$ and β is $k \times 1$.

$$H_0 : \theta = \theta_0; \quad H_1 : \theta \neq \theta_0$$

We have

$$\sqrt{N}(\hat{\theta} - \theta_0) = H\sqrt{N}(\hat{\beta} - \beta_0) \xrightarrow[H_0]{D} N(0, H\Lambda_0 H')$$

Moreover,

$$W_0 = N(\hat{\theta} - \theta_0)(H\Lambda_0 H')^{-1}(\hat{\theta} - \theta_0) \xrightarrow[H_0]{D} \chi_p^2$$

where $\mathbb{E}[\chi_p^2] = p$.

6.3 Long, Short, Auxiliary Regression

$Y \in \mathbb{R}^1$, $X \in \mathbb{R}^K$, $K \in \mathbb{R}^J$. Consider a researcher interested in the conditional distribution of the logarithm of weekly wages ($Y \in \mathbb{R}^1$) given years of completed schooling ($X \in \mathbb{R}^K$) and vector of additional worker attributes. This vector could include variables such as age, childhood test scores, and race. Let W be this $J \times 1$ vector of additional variables.

We can run regression by two ways:

1. Long regression: $\mathbb{E}^*[Y|X, W] = X'\beta_0 + W'\gamma_0$.

2. Short regression: $\mathbb{E}^*[Y|X] = X'b_0$.

Proposition 6.3 (Long Regression)

Long regression is another form of best linear predictor.

$$\begin{aligned}\mathbb{E}^*[Y|X, W] &= \mathbb{E}^*[Y|Z] \\ &= Z' (\mathbb{E}[ZZ'])^{-1} \mathbb{E}[ZY] \\ &= X'\beta_0 + W'\gamma_0\end{aligned}$$

where $\begin{pmatrix} \beta_0 \\ \gamma_0 \end{pmatrix} = \mathbb{E}[ZZ']^{-1} \mathbb{E}[ZY]$, $Z = \begin{pmatrix} X \\ W \end{pmatrix}$.



Proposition 6.4 (Auxiliary Regression)

$$\mathbb{E}^*[W|X] = \Pi_0 X$$

which is multivariate regression. For each row $j = 1, \dots, J$,

$$\mathbb{E}^*[W_j|X] = X'\Pi_{j0}$$

where $\Pi_{j0} = \mathbb{E}[XX']^{-1} \mathbb{E}[XW_j]$ and $\Pi_0 = \begin{pmatrix} \Pi'_{10} \\ \vdots \\ \Pi'_{J0} \end{pmatrix} = \mathbb{E}[WX'] \mathbb{E}[XX']^{-1}$.



Theorem 6.3 (Law of Iterated Linear Predictors (LILP))

$$\mathbb{E}^*[Y|X] = \mathbb{E}^*[\mathbb{E}^*[Y|X, W]|X]$$



Facts: Linear predictor is linear operator, $\mathbb{E}^*[X + Y|W] = \mathbb{E}^*[X|W] + \mathbb{E}^*[Y|W]$.

Let $Y = \mathbb{E}^*[Y|X, W] + u = X'\beta_0 + W'\gamma_0 + u$. Then,

$$\begin{aligned}\mathbb{E}^*[Y|X] &= \mathbb{E}^*[X'\beta_0 + W'\gamma_0 + u|X] \\ &= \mathbb{E}^*[X'\beta_0|X] + \mathbb{E}^*[W'\gamma_0|X] + \mathbb{E}^*[u|X] \\ &= X'\beta_0 + (\Pi_0 X)'\gamma_0 + 0 \\ &= X'(\underbrace{\beta_0 + \Pi'_0 \gamma_0}_{b_0})\end{aligned}$$

Proposition 6.5 (Short Regression)

$$\mathbb{E}^*[Y|X] = X'b_0$$

where $b_0 = \beta_0 + \Pi'_0 \gamma_0$.



6.4 Residual Regression

Let the variation in W unexplained by X .

$$\underbrace{V}_{J \times 1} = \underbrace{W}_{J \times 1} - \underbrace{\mathbb{E}^*[W|X]}_{J \times 1} = W - \Pi_0 X$$

Proposition 6.6 (Residual Regression)

Let $\tilde{Y} = Y - \mathbb{E}^*[Y|X]$,

$$\mathbb{E}^*[\tilde{Y}|V] = V'\gamma_0$$



Proof 6.1

$$Y = X'\beta_0 + W'\gamma_0 + u$$

$$\tilde{Y} = X'\beta_0 - \mathbb{E}^*[Y|X] + W'\gamma_0 + u$$

$$= -X'(\Pi'_0 \gamma_0) + W'\gamma_0 + u$$

$$= V'\gamma_0 + u$$

$$\mathbb{E}^*[\tilde{Y}|V] = V'\gamma_0$$

By long regression,

$$\begin{aligned} \mathbb{E}^*[Y|X, W] &= X'\beta_0 + W'\gamma_0 \\ &= X'b_0 - X'(\Pi'_0 \gamma_0) + W'\gamma_0 \\ &= X'b_0 + V'\gamma_0 \\ &= \mathbb{E}^*[Y|X] + \mathbb{E}^*[\tilde{Y}|V] \end{aligned}$$

Theorem 6.4 (Frisch-Waugh Theorem)

$$\begin{aligned} \mathbb{E}^*[Y|X, V] &= \mathbb{E}^*[Y|X] + \mathbb{E}^*[Y|V] - \mathbb{E}[Y] \\ &= \mathbb{E}^*[Y|X, W] \end{aligned}$$



Lemma 6.2

If $Cov(X, W) = 0$, then

$$\mathbb{E}^*[Y|X, W] = \mathbb{E}^*[Y|X] + \mathbb{E}^*[Y|W] - \mathbb{E}[Y]$$



Proof 6.2

Let $u = Y - \mathbb{E}^*[Y|X, W]$.

$$\begin{aligned} 0 &= \mathbb{E}[uW] \\ &= \mathbb{E}[(Y - \mathbb{E}^*[Y|X] - \mathbb{E}^*[Y|W] + \mathbb{E}[Y])W] \\ &= \underbrace{\mathbb{E}[(Y - \mathbb{E}^*[Y|W])W]}_{=0 \text{ by F.O.C.}} - \underbrace{\mathbb{E}[\mathbb{E}^*[Y|X]]}_{=\mathbb{E}[Y]} \mathbb{E}[W] + \mathbb{E}[Y]\mathbb{E}[W] \end{aligned}$$

6.5 Card-Krueger Model

Consider a model about log-learning based on schooling, ability, luck.

$$Y(s) = \alpha_0 + \beta_0 \underbrace{s}_{\text{schooling } s \in \mathbb{S}} + \underbrace{A}_{\text{ability}} + \underbrace{V}_{\text{luck}}$$

Given a cost function about s :

$$C(s) = \underbrace{C}_{\text{cost heterogeneity}} s + \frac{k_0}{2} s^2$$

Assumption We assume

1. Information set $I_0 = (C, A)$ are known by agent when choosing schooling.
2. V is independent of C, A : $V|C, A \triangleq V$.

Then, the observed schooling s should satisfy

$$\begin{aligned} s &= \arg \max_s \mathbb{E}[Y(s) - C(s) | I_0] \\ &= \arg \max_s \alpha_0 + \beta_0 s + A - Cs - \frac{k_0}{2} s^2 \end{aligned}$$

By F.O.C.

$$\beta_0 - C - k_0 s = 0 \Rightarrow s = \frac{\beta_0 - C}{k_0}$$

1. Long Regression:

$$\mathbb{E}^*[Y|s, A] = \alpha_0 + \beta_0 s + A \quad (\text{LR})$$

2. Short Regression:

$$\mathbb{E}^*[Y|s] = a_0 + b_0 s$$

3. Auxillary Regression:

By the best linear predictor, the $\mathbb{E}^*[A|s]$ can be written as

$$\begin{aligned} \mathbb{E}^*[A|s] &= \mathbb{E}[A] - \frac{\text{Cov}(A, s)}{\text{Var}(s)} \mathbb{E}[s] + \frac{\text{Cov}(A, s)}{\text{Var}(s)} s \\ &= \mathbb{E}[A] - \eta_0 \mathbb{E}[s] + \eta_0 s \end{aligned} \quad (\text{AR})$$

where $\eta_0 = \frac{\text{Cov}(A, s)}{\text{Var}(s)}$ and $s = \frac{\beta_0 - C}{k_0}$ and $\mathbb{E}[s] = \frac{\beta_0 - \mu_C}{k_0}$,

$$\begin{aligned}\text{Cov}(A, s) &= \text{Cov}\left(A, \frac{\beta_0 - C}{k_0}\right) = -\frac{\text{Cov}(A, C)}{k_0} = -\frac{\sigma_{AC}}{k_0} \\ \text{Var}(s) &= \text{Var}\left(\frac{\beta_0 - C}{k_0}\right) = \frac{\sigma_C^2}{k_0^2} \\ \eta_0 &= -k_0 \frac{\sigma_{AC}}{\sigma_C^2} = -k_0 \frac{\sigma_{AC}}{\sigma_A \sigma_C} \frac{\sigma_A}{\sigma_C} = -k_0 \rho_{AC} \frac{\sigma_A}{\sigma_C}\end{aligned}$$

The Auxillary Regression is written as

$$\begin{aligned}\mathbb{E}^*[A|s] &= \mathbb{E}[A] + k_0 \rho_{AC} \frac{\sigma_A}{\sigma_C} \frac{\beta_0 - \mu_C}{k_0} - k_0 \rho_{AC} \frac{\sigma_A}{\sigma_C} s \\ &= \mathbb{E}[A] + \rho_{AC} \frac{\sigma_A}{\sigma_C} (\beta_0 - \mu_C) - k_0 \rho_{AC} \frac{\sigma_A}{\sigma_C} s\end{aligned}\tag{AR-1}$$

Hence, the **Short Regression**

$$\begin{aligned}\mathbb{E}^*[Y|s] &= \mathbb{E}^* [\mathbb{E}^*[Y|s, A]|s] \\ &= \mathbb{E}^* [\alpha_0 + \beta_0 s + A|s] \\ &= \alpha_0 + \beta_0 s + \mathbb{E}^*[A|s] \\ &= \alpha_0 + \underbrace{\mathbb{E}[A] + \rho_{AC} \frac{\sigma_A}{\sigma_C} (\beta_0 - \mu_C)}_{a_0} + \underbrace{\left(\beta_0 - k_0 \rho_{AC} \frac{\sigma_A}{\sigma_C}\right) s}_{b_0}\end{aligned}\tag{SR}$$

6.5.1 Proxy Variable Regression

What if we don't observe A or C . We observe some observed variables W (**proxy variable**) instead.

Assumption *We assume*

1. *Redundancy:* $\mathbb{E}^*[Y|s, A, W] = \mathbb{E}^*[Y|s, A]$ (W doesn't give extra information).
2. *Conditional Uncorrelatedness:* $\mathbb{E}^*[A|s, W] = \mathbb{E}^*[A|W] = \Pi_0 + W'\Pi_W$ (Auxillary Regression).
3. *Conditional Independence:* $C \perp A|W = w$.

The **Proxy Variable Regression** is given by

$$\begin{aligned}\mathbb{E}^*[Y|s, W] &= \mathbb{E}^* [\mathbb{E}^*[Y|s, A, W]|s, W] \\ &= \mathbb{E}^* [\mathbb{E}^*[Y|s, A]|s, W] \\ &= \mathbb{E}^* [\alpha_0 + \beta_0 s + A|s, W] \\ &= \alpha_0 + \beta_0 s + (\Pi_0 + W'\Pi_W) \\ &= (\alpha_0 + \Pi_0) + \beta_0 s + W'\Pi_W\end{aligned}\tag{PVR}$$

A general form of **Proxy Variable Regression** with

1. Long Regression: $\mathbb{E}^*[Y|X, A] = X'\beta_0 + A'\gamma_0$
2. Redundancy: $\mathbb{E}^*[Y|X, A, W] = \mathbb{E}^*[Y|X, A]$

3. Conditional Uncorrelatedness: $\mathbb{E}^*[A|X, W] = \mathbb{E}^*[A|W] = \Pi_0 W$

where Π_0 is $P \times J$, W is $J \times 1$, and A is $P \times 1$.

$$\begin{aligned}\mathbb{E}^*[Y|X, W] &= \mathbb{E}^* [\mathbb{E}^*[Y|X, A, W]|X, W] \\ &= \mathbb{E}^* [\mathbb{E}^*[Y|X, A]|X, W] \\ &= \mathbb{E}^* [X'\beta_0 + A'\gamma_0|X, W] \\ &= X'\beta_0 + \mathbb{E}^*[A|X, W]'\gamma_0 \\ &= X'\beta_0 + W'\Pi_0'\gamma_0\end{aligned}$$

Chapter 7 Machine Learning in Inference

Instead of given a prior distribution of Y , we are given a **training set** $T = (X_i, Y_i)_{i=1}^n$ where i.i.d. $(X_i, Y_i) \sim P$. (Distribution P is unknown).

Risk: $R(\delta) = \mathbb{E}_P [l(\delta(X), Y)]$

The true optimal decision rule is

$$\delta_B = \operatorname{argmin}_{\delta} \mathbb{E}_P [l(\delta(X), Y)]$$

which is can't be computed since we don't know how actually P is.

7.1 Empirical Risk Minimization (ERM)

Instead of computing optimal decision rule with P , we compute the optimal decision rule in the training set:

$$\hat{\delta}_n = \operatorname{argmin}_{\delta} \frac{1}{n} \sum_{i=1}^n l(\delta(X_i), Y_i)$$

The corresponding risk is $R(\hat{\delta}_n) = \mathbb{E}_P [l(\hat{\delta}_n(X), Y)]$. $\Delta R(\hat{\delta}_n) = R(\hat{\delta}_n) - R(\delta) > 0$ always holds.

Consistency: if $\Delta R(\hat{\delta}_n) \rightarrow 0$ as $n \rightarrow \infty$.

7.1.1 Example: Linear MMSE (LMMSE) estimator

Use the decision rule in the class of $\delta(x) = wx$. To find the linear MMSE (LMMSE) estimation $\delta^*(x) = w^*x$:

$$w^* = \operatorname{argmin}_w \mathbb{E}_P [(wx - Y)^2] = \frac{\mathbb{E}[XY]}{\mathbb{E}[X^2]}$$

The rule that minimizes the **empirical risk** is

$$\hat{w} = \operatorname{argmin}_w \frac{1}{n} \sum_{i=1}^n (wx_i - Y_i)^2 = \frac{\frac{1}{n} \sum_{i=1}^n X_i Y_i}{\frac{1}{n} \sum_{i=1}^n X_i^2}$$

The risk of the optimal rule $\delta^* = w^*x$ is $R(\delta^*)$ and the empirical risk under rule $\hat{\delta}(x) = \hat{w}x$ is $R(\hat{\delta}(x))$.

$R(\hat{\delta}) > R(\delta^*)$ always holds, and

$$R(\hat{\delta}) \rightarrow R(\delta^*) \text{ as } n \rightarrow \infty$$

According to CLT:

$$\begin{aligned} \sqrt{n} \left(\frac{1}{n} \sum_i X_i Y_i - \mathbb{E}(XY) \right) &\xrightarrow{d} N(0, \sigma^2) \\ \sqrt{n} \left(\frac{1}{n} \sum_i X_i^2 - \mathbb{E}(X^2) \right) &\xrightarrow{d} N(0, \sigma^2) \end{aligned}$$

Then,

$$\begin{aligned}\frac{1}{n} \sum_i X_i^2 &= \mathbb{E}(X^2) + O\left(\frac{1}{\sqrt{n}}\right) \\ \frac{1}{n} \sum_i X_i Y_i &= \mathbb{E}(XY) + O\left(\frac{1}{\sqrt{n}}\right)\end{aligned}$$

which means the error of estimators

$$\begin{aligned}\hat{w} &= \frac{\mathbb{E}(X^2) + O\left(\frac{1}{\sqrt{n}}\right)}{\mathbb{E}(XY) + O\left(\frac{1}{\sqrt{n}}\right)} = w^* + O\left(\frac{1}{\sqrt{n}}\right) \\ \hat{w} - w^* &= O\left(\frac{1}{\sqrt{n}}\right)\end{aligned}$$

and the error of the risks:

$$\begin{aligned}R(\hat{\delta}) - R(\delta^*) &= \mathbb{E}_P[\hat{w}X - Y]^2 - \mathbb{E}_P[w^*X - Y]^2 \\ &= \mathbb{E}_P[(\hat{w} - w^*)X + w^*X - Y]^2 - \mathbb{E}_P[w^*X - Y]^2 \\ &= \mathbb{E}_P[(\hat{w} - w^*)X]^2 + 2(\hat{w} - w^*)\mathbb{E}_P[X(w^*X - Y)] \\ &= (\hat{w} - w^*)\mathbb{E}_P(X^2) = O\left(\frac{1}{n}\right)\end{aligned}$$

Complexity:

Definition 7.1

A sequence $f(n)$ is $O(1)$ if $\lim_{n \rightarrow \infty} f(n) < \infty$.



Definition 7.2

A sequence $f(n)$ is $O(g(n))$ if $\frac{f(n)}{g(n)}$ is $O(1)$.



Definition 7.3

A sequence $f(n)$ is $o(1)$ if $\lim_{n \rightarrow \infty} \sup f(n) = 0$.



Definition 7.4

A sequence $f(n)$ is $o(g(n))$ if $\lim_{n \rightarrow \infty} \sup \frac{f(n)}{g(n)} = 0$.



Definition 7.5

A sequence $f(n)$ is asymptotic to $g(n)$ if $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 1$. (This is denoted by $f(n) \sim g(n)$ as $a \rightarrow \infty$)



7.1.2 Penalized ERM

$$\delta(x) = \sum_{j=1}^J w_j x^j$$

Pick $J = d$ and use ERM with d dimensional w :

$$\operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n [w^T X_i - Y_i]^2$$

Approach 1: Fix $d \ll n$, use ERM.

Approach 2: (**Penalized ERM**)

$$\min_{\delta} [R_{emp}(\delta) + J(\delta)]$$

$(J(\delta))$ is regularization (penalty) term

7.2 Stochastic Approximation

Robbins and Monro (95)

Problem: Find a root of function $h(x)$. ($f(x) = 0$)

We do not observe $h(x)$ directly, but we observe $Y \sim P_x$ with

- (1) $\mathbb{E}[Y|X=x] = h(x)$
- (2) $(Y|X=x) - h(x)$ is bounded

Example 7.1 $Y = X + Z$ with $\mathbb{E}[Z] = 0$ and Z is bounded.

Assumptions: 1. $h'(x^*) > 0$; 2. x^* is the unique root of h .

SA Algorithm

- Pick Sequence $\{a_n\}$ such that $\sum_{n=1}^{\infty} a_n = \infty$ and $\sum_{n=1}^{\infty} a_n^2 < \infty$ (should converge to 0 but not too quick) e.g. $a_n = n^{-\alpha}$ when $\alpha \in (\frac{1}{2}, 1]$.
- Initialize X_1
- Update for $n = 1, 2, \dots$, $Y_n \sim P(\cdot|X=X_n)$

$$X_{n+1} = X_n - a_n Y_n$$

until convergence.

Theorem 7.1

Under these assumptions

$$X_n \xrightarrow{m.s.} x^* \text{ as } n \rightarrow \infty$$

i.e., $\mathbb{E}(X_n - x^*)^2 \rightarrow 0$ as $n \rightarrow \infty$.



- **Performance Measure** (Convergence rate): the root mean squared error (RMSE) $e_n = \sqrt{\mathbb{E}[(X_n - x^*)^2]}$.
- **Projections:** If x is constrained to live in an interval I , the update rule becomes

$$X_{n+1} = \text{Proj}_x[X_n - a_n Y_n]$$

- **Averaging:**

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i = \frac{X_n}{n} + \bar{X}_{n-1} \frac{n-1}{n}$$

(nicer graph) (The benefits of this smoothing operation are mostly seen in the initial stages of the SA recursion, and do not improve the convergence rate.)

Example 7.2 Let $h(x) = x$, in which case $x^* = 0$. $Y_n = h(X_n) + Z_n = X_n + Z_n$ where noise Z_n is independent of Y_n with $\mathbb{E}[Z_n] = 0$, $Var(Z_n) = 1$ and Z_n is bounded.

Then,

$$\begin{aligned} X_{n+1} &= X_n - a_n(X_n + Z_n) \\ &= (1 - a_n)X_n - a_n Z_n \end{aligned}$$

The MSE,

$$\begin{aligned} e_{n+1}^2 &= \mathbb{E}(X_{n+1} - x^*)^2 = \mathbb{E}(X_{n+1})^2 \\ &= \mathbb{E}[(1 - a_n)X_n - a_n Z_n]^2 \\ &= (1 - a_n)^2 \mathbb{E}X_n^2 + a_n^2 \mathbb{E}Z_n^2 \\ &= (1 - a_n)^2 e_n^2 + a_n^2 \end{aligned}$$

Pick $a_n = n^{-\alpha}$, where $\alpha \in (\frac{1}{2}, 1]$

$$\Rightarrow e_{n+1}^2 = (1 - n^{-\alpha})^2 e_n^2 + n^{-2\alpha}$$

Guess: $e_n = \sqrt{c}n^{-\beta} + H.O.T$

$$c(n+1)^{-2\beta} + H.O.T = (1-n^{-\alpha})^2 cn^{-2\beta} + n^{-2\alpha} + H.O.T$$

(where $(n+1)^{-2\beta} = n^{-2\beta}(1+\frac{1}{n})^{-2\beta} = n^{-2\beta}[1 - 2\beta n^{-1} + O(n^{-2})]$, by Taylor)

$$cn^{-2\beta} - 2c\beta n^{-1-2\beta} + H.O.T = (1-n^{-\alpha})^2 cn^{-2\beta} + n^{-2\alpha} + H.O.T$$

$$-2c\beta n^{-1-2\beta} + H.O.T = -2cn^{-\alpha-2\beta} + n^{-2\alpha} + H.O.T$$

(For $\alpha < 1$), $-2c\beta n^{-1-2\beta}$ is not dominant term.

$$H.O.T = -2cn^{-\alpha-2\beta} + n^{-2\alpha} + H.O.T$$

Identify Power: $2\alpha = \alpha + 2\beta \Rightarrow \beta = \frac{\alpha}{2}$ and $c = \frac{1}{2}$

(For $\alpha = 1$), there are three dominant terms.

$$-2c\beta n^{-1-2\beta} + H.O.T = -2cn^{-1-2\beta} + n^{-2} + H.O.T$$

Identify Power: $2 = 1 + 2\beta \Rightarrow \beta = \frac{1}{2}$ and $-2c\beta = -2c + 1 \Rightarrow c = 1$

$$e_n^2 \sim cn^{-2\beta}$$

To let the convergence rate as fast as possible, we want the β to be as large as possible. Since

$\beta = \frac{\alpha}{2}$, we pick the highest $\alpha = 1 \Rightarrow \beta = \frac{1}{2}, c = 1$.

$$e_n = O(n^{-\frac{1}{2}}) \text{ with } a_n \sim \frac{1}{n}$$

Example 7.3 Let $h(x) = x^3$, in which case $x^* = 0$. $Y_n = h(X_n) + Z_n = X_n^2 + Z_n$ where noise Z_n is independent of Y_n with $\mathbb{E}[Z_n] = 0, Var(Z_n) = 1$ and Z_n is bounded.

Then,

$$\begin{aligned} X_{n+1} &= X_n - a_n(X_n^3 + Z_n) \\ &= X_n - a_nX_n^3 - a_nZ_n \end{aligned}$$

Pick $a_n = n^{-\alpha}, \alpha \in (\frac{1}{2}, 1] \Rightarrow \beta = \frac{1}{6}, \alpha = \frac{2}{3} \Rightarrow e_n \sim O(n^{-\frac{1}{6}})$

7.3 Stochastic Gradient Descent (SGD)

Solve $\min_{x \in \mathbb{R}^n} f(x)$.

We only use a **noisy version** $g(x, z)$ of $f(x)$, where $\mathbb{E}_z[g(x, z)] = f(x)$.

$$\mathbb{E}_z[\nabla_x g(x, z)] = \nabla_x \mathbb{E}_z[g(x, z)] = \nabla f(x)$$

Also pick sequence $\{a_n\}$ such that $\sum_{n=1}^{\infty} a_n = \infty$ and $\sum_{n=1}^{\infty} a_n^2 < \infty$.

SGD

- Initialize X_1
- Update for $n = 1, 2, \dots$,

$$X_{n+1} = X_n - a_n \nabla g(X_n, Z_n)$$

Example 7.4 $f(x) = \frac{1}{2}x^2, x \in \mathbb{R}$. Let Z be a random variable with $\mathbb{E}(Z) = 0, \text{Var}(Z) = 1$.

$$\begin{aligned} g(x, Z) &= \frac{1}{2}(x + Z)^2 - \frac{1}{2} \\ \mathbb{E}[g(x, Z)] &= \frac{1}{2}x^2 = f(x) \\ \nabla_x g(x, Z) &= x + Z \Rightarrow \mathbb{E}[\nabla_x g(x, Z)] = \nabla f(x) \\ X_{n+1} &= X_n - a_n(X_n + Z_n) \end{aligned}$$

which is the same as the stochastic approximation.

Main Results: (Suppose the unique minimum is x^*)

- (1) Convergence: $e_n \rightarrow 0$ as $n \rightarrow \infty$.
- (2) Convergence Rate: To achieve $\mathbb{E}[f(X_n)] - f(x^*) < \varepsilon$, we need $n = O(\frac{1}{\varepsilon})$ if f is twice continuously differentiable and strongly convex.

GD has linear convergence $\Rightarrow e_n = O(e^{-cn})$; Solve $\varepsilon = O(e^{-cn}) \Rightarrow n = O(\ln \frac{1}{\varepsilon})$. (**SGD is much worse than GD**, cost more.)

7.4 SGD Application to Empirical Risk Minimization (ERM)

ERM problem is

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n L(\delta_w(X_i), Y_i)$$

$R_{emp}(w) = \frac{1}{n} \sum_{i=1}^n L(\delta_w(X_i), Y_i)$ is the empirical risk (e.g. $\delta_w(x) = w^T x, L(\hat{y}, y) = (\hat{y} - y)^2$) To make w more visible, we can write

$$R_{emp}(w) = \frac{1}{n} \sum_{i=1}^n L(\delta_w(X_i), Y_i) = \frac{1}{n} \sum_{i=1}^n Q(X_i, Y_i, w)$$

For penalized ERM we would similarly have

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n Q(X_i, Y_i, w) + J(w)$$

$J(w)$ is the penalty (regularization) term.

In the problem $\min_{W \in \mathbb{R}^n} R_{emp}(w) = \frac{1}{n} \sum_{i=1}^n Q(X_i, Y_i, w)$

7.4.1 Different Gradient Descent for ERM

GD • Initialize W_1

- Update for $k \geq 1$, Update: $W_{k+1} = W_k - a_k \frac{1}{n} \sum_{i=1}^n \nabla Q(X_i, Y_i, W_k)$

Computational cost: The computational cost of GD is $O(dn)$ operations per iteration. Since GD has exponential convergence, the number of iterations needed to reach an optimization error of ρ is $O(\log \frac{1}{\rho})$. Hence GD incurs a total computational cost of $O(dn \log \frac{1}{\rho})$ to reach a solution W_k such that $R_{emp}(W_k) \leq \min_W R_{emp}(W) + \rho$

SGD • Initialize W_1

- Update for $k \geq 1$,

Step 1: Pick i uniformly over $\{1, \dots, n\}$

Step 2: $W_{k+1} = W_k - a_k \nabla Q(X_i, Y_i, W_k)$

Computational cost: After k iterations, $\mathbb{E}[R_{emp}(W_k)] \leq \min_W R_{emp}(W) + \rho$, for $k = O(\frac{1}{\rho})$ and f twice differentiable and strongly convex. The cost per iteration is $O(d)$ (independent of n), so the total computational cost is $O(\frac{d}{\varepsilon})$.

7.4.2 Constraints on Learning Problem

Why achieving a low value of ρ is useful (low error of $R_{emp}(\cdot)$), since the cost function $R_{emp}(\cdot)$ is only a surrogate for the actual risk $R(\cdot)$?

Typically, $d = O(n^b)$ where $0 < b < 1$.

For any numerical algorithm producing a decision rule $\tilde{\delta}_n$, the excess risk (compared to Bayes rule δ_B) can be expressed as the sum of three terms:

$$\begin{aligned}\Delta R(\tilde{\delta}_n) &= R(\tilde{\delta}_n) - R(\delta_B) \\ &= [R(\tilde{\delta}_n) - R(\hat{\delta}_n)] + [R(\hat{\delta}_n) - R(\delta^*)] + [R(\delta^*) - R(\delta_B)]\end{aligned}$$

where

$$\delta_B = \text{Bayes rule}$$

$$\delta^* = \text{best rule in } D = \operatorname{argmin}_{\delta \in D} R(\delta)$$

$$\hat{\delta}_n = \operatorname{argmin}_{\delta \in D} R_{emp}(\delta)$$

$$\tilde{\delta}_n = \text{solution of the algorithm after } k \text{ iterations}$$

(Note: D is the set of all available decision rule in approximation (e.g. all linear parameters $\{W, b\}$), which can't be better than Bayes rule)

The expected excess risk

$$\epsilon = \mathbb{E}[\Delta R(\tilde{\delta}_n)] = \underbrace{\mathbb{E}[R(\tilde{\delta}_n) - R(\hat{\delta}_n)]}_{\text{Comp. Error}=\rho} + \underbrace{\mathbb{E}[R(\hat{\delta}_n) - R(\delta^*)]}_{\text{Est. Error}=O(\frac{d}{n})} + \underbrace{\mathbb{E}[R(\delta^*) - R(\delta_B)]}_{\text{Approx. Error}=O(d^{-\beta})}$$

Estimation error increases as d increases, but *approximation error* decreases as d increases. To minimize the excess risk, we want to balance the last two items, that is $O(\frac{d}{n}) = O(d^{-\beta})$: solve

$$\begin{aligned} \frac{d}{n} = d^{-\beta} &\Rightarrow d^{1+\beta} = n \Rightarrow d = n^{\frac{1}{1+\beta}} \\ \Rightarrow \text{the last two items } O\left(\frac{d}{n}\right) &= O(d^{-\beta}) = O(n^{-\gamma}) \end{aligned}$$

where $\gamma = \frac{\beta}{1+\beta} \in (0, 1]$ is a constant.

To balance the three items, we want

$$\rho = O(n^{-\gamma}) \Rightarrow n = O(\rho^{-\frac{1}{\gamma}}) \text{ and } d = O(n^{\frac{1}{1+\beta}})$$

The update rule $W_{k+1} = W_k - a_k \nabla Q(X_k, Y_k, W_k)$, where $i \sim \text{Uniform}\{1, 2, \dots, n\}$

Relation to Online Learning: When the training data are made available sequentially (instead of in a batch as assumed here), online learning can be used to sequentially learn the decision rules (or the weights that parameterize the decision rule).

Variations on Basic SGD: mini batch: replace S by a subset B and n by $|B|$

$$\frac{1}{|B|} \sum_{i \in B} Q(X_i, Y_i, W_k)$$

Averaging SGD:

$$\bar{W}_n = \frac{1}{n} \sum_{i=1}^n W_i = \frac{W_n}{n} + \bar{W}_{n-1} \frac{n-1}{n}$$

SVRG (Stochastic Variance Randomized Gradient): R. Johnson and T. Zhang, “Accelerating Stochastic Gradient Descent using Predictive Variance Reduction,” *Proc. NIPS* 2013.

Unsupervised learning: If no explanatory variable X is present, the problem reduces to

$$\min_w \frac{1}{n} \sum_{i=1}^n Q(Y_i, w)$$

which finds applications to various unsupervised learning problems. For instance the k -means clustering algorithm partitions n data points $y_i, 1 \leq i \leq n$ in \mathbb{R}^d into k clusters with centroids $w_j, 1 \leq j \leq k$, in a way that minimizes the within-cluster sum-of-squares: $\text{WCSS} = \frac{1}{n} \sum_{i=1}^n \min_{1 \leq j \leq k} \|y_i - w_j\|^2$. Using the formalism of (7), we have $Q(y, w) = \min_{1 \leq j \leq k} \|y - w_j\|^2$ where $y \in \mathbb{R}^d$ and $w = \{w_j\}_{j=1}^k \in \mathbb{R}^{d \times k}$ is a matrix whose k columns are the centroid vectors.

Chapter 8 Stochastic Integration Methods

Integral $I = \int_X f(x)dx$.

8.1 Deterministic Methods (Better in Low Dimension)

8.1.1 Riemann Integration

Riemann integral: approximation integral I of $f(x)$ in $[a, b]$ with

$$\hat{I}_n = \sum_{i=1}^n \underbrace{(x_i - x_{i-1})}_{\frac{b-a}{n}} f(x_i)$$

where $x_i = a + \frac{b-a}{n}i$. We can also denote $\hat{f}(x) = f(x_i)$ if $x \in (x_{i-1}, x_i]$

The error $|\hat{I}_n - I| = \int_a^b |\hat{f}(x) - f(x)|dx$

Assume f is differentiable and $\max_x |f'(x)| = c < \infty$, then $|\hat{f}(x) - f(x)| \leq \frac{b-a}{n}c$

$$\Rightarrow |\hat{I}_n - I| \leq \int_a^b \frac{b-a}{n}c dx = \frac{(b-a)^2}{n}c$$

That is $n \sim O(\varepsilon^{-1})$

8.1.2 Trapezoidal Rule

Using average can be better.

$$\hat{I}_n = \sum_{i=1}^n \underbrace{(x_i - x_{i-1})}_{\frac{b-a}{n}} \frac{f(x_i) + f(x_{i-1})}{2}$$

The upper bound of error is $|\hat{I}_n - I| \leq \frac{C}{n^2}$ for some constant C .

That is $n \sim O(\varepsilon^{-\frac{1}{2}})$

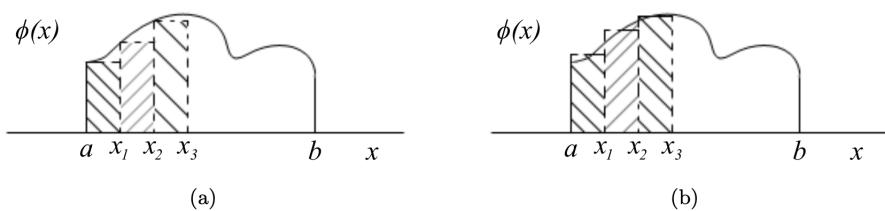


Figure 8.1: (a) Riemann approximation; (b) Trapezoidal approximation.

8.1.3 Multidimensional Integration

When we want to do integral in high dimension, it will be really hard.

For d -dimensional integrals, the trapezoidal rule yields an approximation error $|\hat{I}_n - I| \leq \frac{C}{n^{\frac{d}{2}}}$ for some constant C . That is $n \sim O\left(\varepsilon^{-\frac{d}{2}}\right)$. n needs to increase exponentially with d to achieve a target approximation error ε .

This phenomenon is known as the curse of dimensionality.

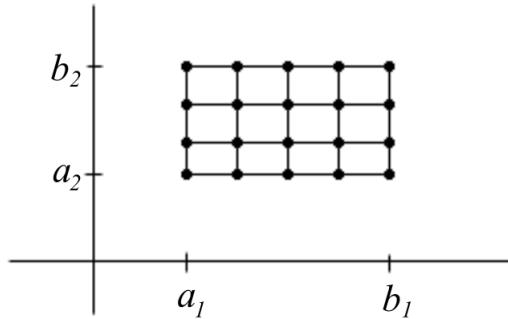


Figure 8.2: Two-dimensional integration using regular grid.

8.2 Stochastic Methods (Better in High Dimension)

8.2.1 Classical Monte Carlo Integration

Compute the expectation

$$\xi = \mathbb{E}_p[h(x)] = \int_X \underbrace{p(x)h(x)}_{f(x)} dx$$

The methods described below can be used to solve the following problems: (1) General $\int_X f$; (2) Compute the probability of falling into a subset $a \subset X : P(a) = \int_a p(x)dx$, where $h(x) = \mathbf{1}_{x \in a}$

The Monte Carlo approach is as follows: Given X_1, X_2, \dots, X_n drawn i.i.d from the pdf p , estimate ξ by the empirical average

$$\hat{\xi}_n = \frac{1}{n} \sum_{i=1}^n h(X_i)$$

$$\mathbb{E}_p[\hat{\xi}_n] = \mathbb{E}_p[h(X)] = \xi. \quad \hat{\xi}_n \xrightarrow{a.s.} \xi \text{ as } n \rightarrow \infty \text{ by SLLW.}$$

$$Var(\hat{\xi}_n - \xi) = Var(\hat{\xi}_n) = \frac{1}{n} Var[h(x)] = O\left(\frac{1}{n}\right) \Rightarrow sd(\hat{\xi}_n) = \frac{\sqrt{Var[h(x)]}}{\sqrt{n}}$$

$$\text{That is } n \sim O\left(n^{-\frac{1}{2}}\right)$$

The stochastic methods **outperform** when the deterministic ones for dimensions $d > 4$ and are **worse** for $d < 4$.

8.2.2 Importance Sampling

Draw $X_i, i = 1, \dots, n$ i.i.d from pdf q

$$\hat{\xi}_n = \frac{1}{n} \sum_{i=1}^n \frac{p(X_i)}{q(X_i)} h(X_i)$$

It is an unbiased estimator of ξ

$$\mathbb{E}[\hat{\xi}_n] = \mathbb{E}_q\left[\frac{p(X_i)}{q(X_i)} h(X_i)\right] = \int_X p(x)h(x)dx = \xi$$

$\hat{\xi}_n \xrightarrow{a.s.} \xi$ as $n \rightarrow \infty$ by SLLW.

Its variance is

$$\begin{aligned} Var_q(\hat{\xi}_n) &= \frac{1}{n} Var_q \left[\frac{p(X_i)}{q(X_i)} h(X_i) \right] \\ &= \frac{1}{n} \left(\int_X \frac{p^2(x)}{q(x)} h^2(x) dx - \xi^2 \right) \end{aligned}$$

The idea of importance sampling is to find a good q such that

$$Var_q(\hat{\xi}_n) < Var_p(\hat{\xi}_n)$$

Error Measure

The *relative error* of the importance-sampling estimator is defined as

$$\delta_{\text{rel}}(\hat{\xi}_n) \triangleq \frac{\sqrt{\text{Var}_q(\hat{\xi}_n)}}{\xi} = \sqrt{\frac{\text{Var}_q\left[\frac{p(X)}{q(X)} h(X)\right]}{\xi^2 n}}.$$

The *number* of simulations needed to achieve a relative error of δ is

$$n_{IS}(\delta) = \frac{\text{Var}_q\left[\frac{p(X)}{q(X)} h(X)\right]}{\xi^2 \delta^2}.$$

The *gain relative to a Monte Carlo simulation* is defined as

$$\Gamma = \frac{n_{MC}(\delta)}{n_{IS}(\delta)} = \frac{\text{Var}_p[h(X)]}{\text{Var}_q\left[\frac{p(X)}{q(X)} h(X)\right]}$$

In the example of $\xi = P(a)$, $h(x) = \mathbf{1}_{x \in a}$: Suppose $\xi = P(a) \approx 10^{-9}$ (small), $\hat{\xi}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i \in a}$. $\mathbf{1}_{X_i \in a}$ is *Bernoulli*(ξ). We have $Var(\hat{\xi}_n) = \frac{\xi(1-\xi)}{n}$.

We can use relative error to measure

$$\delta_{\text{rel}}(\hat{\xi}_n) = \frac{\sqrt{Var(\hat{\xi}_n)}}{\xi} = \sqrt{\frac{1-\xi}{n\xi}}$$

and the number of simulation need to get relative error δ is

$$n_{IS}(\delta) = \frac{1-\delta}{\xi \delta^2}.$$

Find the optimal q :

$$\min_q \int_X \frac{p^2(x)}{q(x)} h^2(x) dx - \xi^2$$

write

$$\int_X \frac{p^2(x)}{q(x)} h^2(x) dx - \xi^2 = \mathbb{E}_q \left[\left(\underbrace{\frac{p(x)}{q(x)} h(x)}_Z \right)^2 \right]$$

Since x^2 is convex function, by Jensen's inequality

$$\mathbb{E}_q \left[\left(\underbrace{\frac{p(x)}{q(x)} h(x)}_Z \right)^2 \right] \geq \left(\mathbb{E}_q \left[\underbrace{\frac{p(x)}{q(x)} h(x)}_Z \right] \right)^2$$

This equality holds if and only if $\frac{p(x)}{q(x)} h(x) = \alpha, \forall x \in X$, α is a constant.

Since q is pdf., we can infer

$$q(x) = \frac{p(x)h(x)}{\int_X p(x)h(x)dx}$$

which is as hard as the original problem. In practice, one is content to find a “good” q that assigns high probability to the important region where $p(x)h(x)$ is large. Ideally the ratio $\frac{p(x)}{q(x)} h(x)$ would be roughly constant over X .

Chapter 9 Particle Filtering

Kalman filtering is used in tracking problems (dynamic models). Particle Filtering is an extension of Kalman filtering.

9.1 Kalman Filtering (Linear Dynamic System)

1. Unknown state sequence $X_t \in \mathbb{R}^m, t = 0, 1, 2, \dots$
2. Observations $Y_t \in \mathbb{R}^k, t = 0, 1, 2, \dots$
3. $X_{t+1} = F_t X_t + U_t, F_t \in \mathbb{R}^{m \times m}, U_t \sim P_{U_t}$
4. $Y_t = H_t X_t + V_t, H_t \in \mathbb{R}^{k \times m}, V_t \sim P_{V_t}$

We want to solve two problems

1. **Estimation Problem:** Evaluate Linear MMSE (LMMSE) of X_t given $Y_{0:t}$.

$$\hat{X}_{t|t} = W Y_{0:t} + b$$

2. **Prediction Problem:** Predict Linear MMSE (LMMSE) of $X_{t+1|t}$ given $Y_{0:t}$. (Really hard)

We can solve closed-form solutions.

9.2 Particle Filtering (Nonlinear Dynamic System)

Particle filtering is a nonlinear form of Kalman filtering, which doesn't have closed-form solutions.

We consider a Nonlinear Dynamic System

$$X_{t+1} \sim q(\cdot | X_t)$$

$$Y_t \sim r(\cdot | X_t)$$

$$t = 0, 1, 2, \dots$$

where $q(X_{t+1}|X_t)$ is the transition probability distribution, and $r(Y_t|X_t)$ is the conditional probability distribution for the observations. Hence, X_t is a Markov process and Y_t follows a Hidden Markov Model (HMM).

We also consider these two problems.

1. **Estimation Problem:** Evaluate X_t given $Y_{0:t}$.
2. **Prediction Problem:** Predict $X_{t+1|t}$ given $Y_{0:t}$.

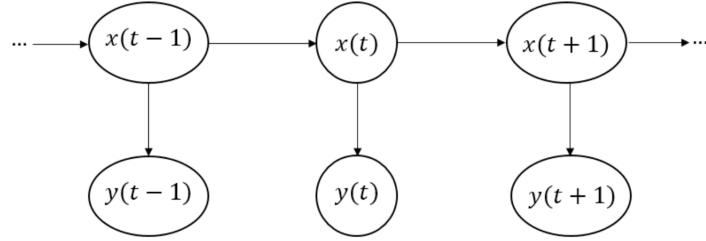


Figure 9.1: Hidden Markov Model

9.2.1 Bayesian Recursive Filtering

In this section we use Bayesian approach and use MMSE estimation $l(\hat{x}_t, x_t) = \|x_t - \hat{x}_t\|^2$

Estimation and prediction in conditional forms are:

$$\begin{aligned}\hat{X}_{t|t} &= \mathbb{E}[X_t|Y_{0:t}] = \int_{\mathbb{R}^m} x_t P(X_t|Y_{0:t}) dx_t \\ \hat{X}_{t+1|t} &= \mathbb{E}[X_{t+1}|Y_{0:t}] = \int_{\mathbb{R}^m} x_{t+1} P(X_{t+1}|Y_{0:t}) dx_{t+1}\end{aligned}$$

Apparently the posterior p.d.f cannot be evaluated due to the curse of dimensionality as t increases. However, they can in principle be evaluated *recursively* using the following two-step procedure.

Step 1: Prediction. $P(X_{t+1}|Y_{0:t})$ can be expressed in term of $P(X_t|Y_{0:t})$:

$$\begin{aligned}P(X_{t+1}|Y_{0:t}) &= \int_{\mathbb{R}^m} P(X_{t+1}, X_t|Y_{0:t}) dx_t \\ &= \int_{\mathbb{R}^m} P(X_{t+1}|X_t, Y_{0:t}) P(X_t|Y_{0:t}) dx_t \\ &= \int_{\mathbb{R}^m} q(X_{t+1}|X_t) P(X_t|Y_{0:t}) dx_t\end{aligned}$$

Step 2: Update. We can also express $P(X_t|Y_{0:t})$ in terms of $P(X_t|Y_{0:t-1})$

$$\begin{aligned}P(X_t|Y_{0:t}) &= P(X_t|Y_t, Y_{0:t-1}) \\ &= \frac{P(Y_t|X_t, Y_{0:t-1}) P(X_t|Y_{0:t-1})}{P(Y_t|Y_{0:t-1})} \\ &= \frac{r(Y_t|X_t) P(X_t|Y_{0:t-1})}{\int_{\mathbb{R}^m} r(Y_t|X_t) P(X_t|Y_{0:t-1}) dx_t}\end{aligned}$$

9.2.2 Particle Filter (bootstrap filter)

Suppose we have n i.i.d. samples of X_t drawn from $p(x_t|Y_{0:t})$: $X_t(1), X_t(2), \dots, X_t(n)$.

$$X_t(i) \sim p(\cdot|Y_{0:t}), 1 \leq i \leq n \quad (\text{Sample 1})$$

We can use above recursive filtering method to generate estimation of X_{t+1} .

Step 1: Prediction. Using the transition probability $q(\cdot|X_t(i))$, $1 \leq i \leq n$ to generate n independent random variables

$$X_{t+1}^*(i) \sim q(\cdot|X_t(i)), 1 \leq i \leq n \quad (\text{Sample 2})$$

Step 2: Update. Upon receiving a new measurement y_{t+1} , evaluate the *importance weights* (nonnegative and summing to 1)

$$w_i = \frac{r(y_{t+1}|X_{t+1}^*(i))}{\sum_{j=1}^n r(y_{t+1}|X_{t+1}^*(j))}, \quad 1 \leq i \leq n$$

Then we resample n times from the set $\{X_{t+1}^*(i)\}_{i=1}^n$ with respective probabilities $\{w_i\}_{i=1}^n$, obtaining i.i.d samples $\{X_{t+1}(j)\}_{j=1}^n$ with probabilities

$$Pr[X_{t+1}(j) = X_{t+1}^*(i)] = w_i, \quad 1 \leq i, j \leq n \quad (\text{Sample 3})$$

By the **weighted bootstrap theorem**, as $n \rightarrow \infty$, the distribution of the resampled $\{X_{t+1}(j)\}_{j=1}^n$ converges to the desired posterior.

Potential issues: 1. n is not large enough. 2. Sample impoverishment

Chapter 10 EM Algorithm

The ML estimator: $\hat{\theta}_{ML} = \operatorname{argmax}_{\theta \in S} \ln p_{\theta}(y)$. Numerical evaluation of maximum-likelihood (ML) estimates is often difficult. The likelihood function may have multiple extreme and the parameter θ may be multidimensional, all of which are problematic for any numerical algorithm.

Maximum-Likelihood (ML) Estimation

Given a vector \vec{y} , find the θ that maximizes $p_{\theta}(\vec{y}) = \prod_{i=1}^n P(y_i|\theta)$

$$\hat{\theta}_{ML} = \operatorname{argmax}_{\theta \in S} \ln p_{\theta}(\vec{y}) = \operatorname{argmax}_{\theta \in S} \sum_{i=1}^n \ln P(y_i|\theta)$$

Solving the closed-form solution is quite hard sometimes, so we may use EM algorithm.

10.1 General Structure of the EM Algorithm

What we want to estimate:

$\theta \in S$ is an unknown parameter that we want to estimate.

What we know:

To help us solve the solution, we construct an unobservable vector \vec{z} corresponding to \vec{y} .

1. There is a complete data space Z and an incomplete data space Y .
2. The reality is $z \in Z$, which has p.d.f $P(z|\theta)$. ($\ln P(z|\theta)$'s derivative should be constructed to be easy.)
3. **Instead of observing the z directly, we can observe $y = h(z) \in Y$ which has p.d.f $P(y|\theta)$.**
4. $h(z) = y$ is a many-to-one mapping.

$$P(y, z|\theta) = P(z|\theta), \quad \forall z \in h^{-1}(y)$$

5. We can infer that the relationship between $P(z|\theta)$ and $P(y|\theta)$ is

$$P(z|\theta) = P(z|y, \theta)P(y|\theta), \quad \forall z \in h^{-1}(y)$$

$$P(y|\theta) = \sum_{z \in h^{-1}(y)} P(z|\theta), \quad \forall y$$

6. For any function f , $\mathbb{E}_z[f(z)|y]$ depends on the p.d.f.

$$\mathbb{E}_{z|\theta}[f(z)|y] = \sum_{z \in h^{-1}(y)} P(z|y, \theta)f(z)$$

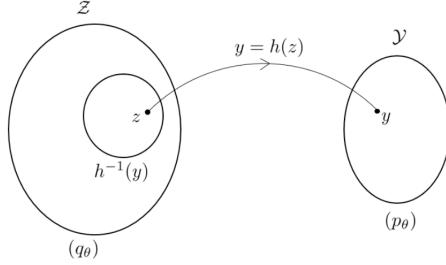


Figure 10.1: Complete and incomplete data spaces Z and Y .

EM Algorithm

Instead of computing the $P(\vec{y}|\theta)$ directly, we use the relationship $h(z) = y$ and $P(\vec{z}|\theta)$ to estimate θ .

Suppose we have a prior belief of the relationship between y and z : $P(z|y, \theta^{(k)})$. Given \vec{y} , since maximizing $\ln P(\vec{y}|\theta)$ is hard, we maximize the expected value of $\ln P(\vec{z}|\theta)|\vec{y}$ under the prior belief (i.e., finding the θ that can properly represent the relationship between \vec{y} and \vec{z}), that is

$$\begin{aligned}\theta &= \operatorname{argmax}_{\theta} \mathbb{E}_{z|y, \theta^{(k)}} [\ln P(\vec{z}|\theta)|\vec{y}] \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^n \sum_{z_i \in h^{-1}(y_i)} P(z_i|y_i, \theta^{(k)}) \ln P(z_i|\theta)\end{aligned}$$

EM algorithm alternates between Expectation (E) and Maximization (M) steps:

1. Initialize $\hat{\theta}^{(0)}$
2. For $k = 0, 1, 2, \dots$

Expectation (E)-Step: Compute

$$\begin{aligned}Q(\theta|\hat{\theta}^{(k)}) &= \mathbb{E}_{z|y, \hat{\theta}^{(k)}} [\ln P(\vec{z}|\theta)|\vec{y}] \\ &= \sum_{i=1}^n \sum_{z_i \in h^{-1}(y_i)} P(z_i|y_i, \theta^{(k)}) \ln P(z_i|\theta)\end{aligned}$$

Maximization (M)-Step

$$\hat{\theta}^{(k+1)} = \operatorname{argmax}_{\theta \in S} Q(\theta|\hat{\theta}^{(k)})$$

Definition 10.1

θ^* is a stable point of the EM algorithm if \exists subsequence that converges to θ^* .

e.g. $1, 3, \frac{1}{2}, 3, \frac{1}{3}, 3, \dots \frac{1}{n}, 3, \dots$



10.2 Example 1: Variance Estimation

Observation $Y = S + N$, $S \sim \mathcal{N}(0, \theta)$ is independent of $N \sim \mathcal{N}(0, \theta) \Rightarrow Y \sim \mathcal{N}(0, \theta + 1)$. $p_\theta(y) = \frac{1}{\sqrt{2\pi(\theta+1)}} e^{-\frac{y^2}{2(\theta+1)}}$. We want to estimate θ .

10.2.1 Maximum-Likelihood (ML) Estimation

$$\ln p_\theta(y) = -\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\theta + 1) - \frac{y^2}{2(\theta + 1)}$$

take derivation of θ to be equal to 0

$$-\frac{1}{2(\theta + 1)} + \frac{y^2}{2(\theta + 1)^2} = 0$$

We can get

$$\hat{\theta} = y^2 - 1$$

Then,

$$\hat{\theta}_{ML} = \begin{cases} 0, & y^2 \leq 1 \\ y^2 - 1, & y^2 > 1 \end{cases}$$

10.2.2 EM Algorithm

Let $Z = (S, N)$, $y = h(z) = s + n$.

$$q_\theta(z) = q_\theta(s, n) = \frac{1}{\sqrt{2\pi\theta}} e^{-\frac{s^2}{2\theta}} \frac{1}{\sqrt{2\pi}} e^{-\frac{n^2}{2}}$$

Then

$$\ln q_\theta(z) = \ln \frac{1}{\sqrt{2\pi}} e^{-\frac{n^2}{2}} - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\theta) - \frac{s^2}{2\theta}$$

E-Step: Compute

$$\begin{aligned} Q(\theta|\hat{\theta}^{(k)}) &= \mathbb{E}_{z|\hat{\theta}^{(k)}} [\ln q_\theta(z)|Y = y] \\ &= \sum_{z \in h^{-1}(y)} q_{\hat{\theta}^{(k)}}(z) \ln q_\theta(z) \\ &= \ln \frac{1}{\sqrt{2\pi}} e^{-\frac{n^2}{2}} - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\theta) - \frac{\mathbb{E}_{z|\hat{\theta}^{(k)}}(s^2)}{2\theta} \end{aligned}$$

M-Step

$$\begin{aligned}\hat{\theta}^{(k+1)} &= \operatorname{argmax}_{\theta \in S} Q(\theta | \hat{\theta}^{(k)}) \\ 0 &= -\frac{1}{2\hat{\theta}^{(k+1)}} + \frac{\mathbb{E}_{z|\hat{\theta}^{(k)}}(s^2)}{2(\hat{\theta}^{(k+1)})^2} \\ \hat{\theta}^{(k+1)} &= \mathbb{E}_{z|\hat{\theta}^{(k)}}(s^2) = \frac{\hat{\theta}^{(k)}}{\hat{\theta}^{(k)} + 1} \left(\frac{\hat{\theta}^{(k)}}{\hat{\theta}^{(k)} + 1} y^2 + 1 \right)\end{aligned}$$

Then we can solve the stable point

$$\begin{aligned}\hat{\theta}^* &= \frac{\hat{\theta}^*}{\hat{\theta}^* + 1} \left(\frac{\hat{\theta}^*}{\hat{\theta}^* + 1} y^2 + 1 \right) \\ \Rightarrow \hat{\theta}^* &= 0, \hat{\theta}^* = y^2 - 1\end{aligned}$$

According to the relation between $\hat{\theta}^{(k)}$ and $\hat{\theta}^{(k+1)}$, we can infer

$$\hat{\theta}^* = \begin{cases} 0, & y^2 \leq 1 \\ y^2 - 1, & y^2 > 1 \end{cases}$$

10.3 Example 2: Estimation of Gaussian Mixtures

Assume the data $\vec{Y} = \{Y_i, 1 \leq i \leq n\} \in \mathbb{R}^n$, are drawn iid from a pdf $p_\theta(y)$ which is the mixture of m univariate Gaussians with respective probabilities $\pi(j)$, means μ_j , and variances σ_j^2 , for $1 \leq j \leq m$:

$$\begin{aligned}p_\theta(y|j) &= \phi(y; \mu_j, \sigma_j^2) \\ p_\theta(y) &= \sum_{j=1}^m \pi(j) \phi(y; \mu_j, \sigma_j^2), \quad y \in \mathbb{R}\end{aligned}$$

where

$$\phi(y; \mu, \sigma^2) \triangleq \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y-\mu)^2}{2\sigma^2} \right\}$$

denotes the Gaussian pdf with mean μ and variance σ^2 .

10.3.1 Unknown Means: ML estimation is hard

We initially assume that $\{\pi(j)\}$ and $\{\sigma_j^2\}$ are given and that we only need to estimate the means $\{\mu_j\}$. Thus, $\theta = \mu \in \mathbb{R}^m$.

Unfortunately the ML estimator cannot be derived in closed form. Indeed, the loglikelihood function for θ is

$$\ln \prod_{i=1}^n p_\theta(y_i) = \sum_{i=1}^n \ln p_\theta(y_i) = \sum_{i=1}^n \ln \sum_{j=1}^m \pi(j) \phi(y_i; \mu_j, \sigma_j^2)$$

and maximizing it is a m -dimensional, nonconcave maximization problem.

Taking the derivative of $\mu_j, j = 1, \dots, m$,

$$\begin{aligned} 0 &= \frac{1}{\sigma_j^2} \sum_{i=1}^n (y_i - \mu_j) \frac{\pi(j)\phi(y_i; \mu_j, \sigma_j^2)}{\sum_{j=1}^m \pi(j)\phi(y_i; \mu_j, \sigma_j^2)} \\ &= \frac{1}{\sigma_j^2} \sum_{i=1}^n (y_i - \mu_j) \pi_\theta(j|y_i) \end{aligned}$$

where $\pi_\theta(j|y_i) \triangleq \frac{\pi(j)\phi(y_i; \mu_j, \sigma_j^2)}{\sum_{j=1}^m \pi(j)\phi(y_i; \mu_j, \sigma_j^2)}$. The system may have multiple solutions corresponding to local maxima or even local minima or saddle points of the likelihood function.

10.3.2 Unknown Means: EM Algorithm

There is a complete data $Z_i = (J_i, Y_i), i = 1, \dots, n$, where J_i is the random label that was drawn to produce Y_i .

$z = \{j_i, y_i\}_{i=1}^n$ is the sample.

$$\begin{aligned} q_\theta(z) &= \prod_{i=1}^n (\pi(j_i)p_\theta(y_i|j_i)) \\ \ln q_\theta(z) &= \sum_{i=1}^n [\ln \pi(j_i) + \ln p_\theta(y_i|j_i)] \end{aligned}$$

Initialize $\hat{\theta}^{(0)}$

Iteration:

$$\begin{aligned} Q(\theta|\hat{\theta}^{(k)}) &= \sum_{i=1}^n \mathbb{E}_{\hat{\theta}^{(k)}} [\ln \pi(j_i) + \ln p_\theta(y_i|j_i)|Y_i = y_i] \\ &= \sum_{i=1}^n \sum_{j=1}^m \pi_{\hat{\theta}^{(k)}}(j|y_i) [\ln \pi(j) + \ln p_\theta(y_i|j)] \\ &= cst - \sum_{i=1}^n \sum_{j=1}^m \pi_{\hat{\theta}^{(k)}}(j|y_i) \frac{(y_i - \mu_j)^2}{2\sigma_j^2} \\ &= cst - \sum_{i=1}^n \sum_{j=1}^m \frac{\pi(j)\phi(y_i; \hat{\mu}_j^{(k)}, \sigma_j^2)}{\sum_{j=1}^m \pi(j)\phi(y_i; \hat{\mu}_j^{(k)}, \sigma_j^2)} \frac{(y_i - \mu_j)^2}{2\sigma_j^2} \end{aligned}$$

where $\ln p_\theta(y_i|j) = -\frac{1}{2} \ln(2\pi\sigma_j^2) - \frac{(y_i - \mu_j)^2}{2\sigma_j^2}$.

Take derivative of μ_j ,

$$\begin{aligned} 0 &= \frac{\partial Q(\theta|\hat{\theta}^{(k)})}{\partial \mu_j} = \sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j|y_i) \frac{(y_i - \mu_j)}{\sigma_j^2} \\ \hat{\mu}_j^{(k+1)} &= \frac{\sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j|y_i)y_i}{\sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j|y_i)} \end{aligned}$$

Recall the $\hat{\theta}_{ML}$

$$\hat{\theta}_{ML,j} = \frac{\sum_{i=1}^n \pi_{\hat{\theta}_{ML}}(j|y_i)y_i}{\sum_{i=1}^n \pi_{\hat{\theta}_{ML}}(j|y_i)}$$

$\hat{\theta}_{ML}$ is the stable point. (if exist)

10.3.3 Unknown Mixture Probabilities, Means and Variances

ML Estimation:

If $\theta \triangleq \{\pi(j), \mu_j, \sigma_j^2, 1 \leq j \leq m\}$ is unknown, the ML estimator $\hat{\theta}_{\text{ML}}$ satisfies the following nonlinear system of equations:

$$\begin{aligned}\hat{\mu}_{\text{ML},j} &= \frac{\sum_{i=1}^n y_i \pi_{\hat{\theta}^{(k)}}(j | y_i)}{\sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j | y_i)} \\ \hat{\sigma}_{\text{ML},j}^2 &= \frac{\sum_{i=1}^n (y_i - \hat{\mu}_{\text{ML},j})^2 \pi_{\hat{\theta}^{(k)}}(j | y_i)}{\sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j | y_i)} \\ \hat{\pi}_{\text{ML}}(j) &= \frac{1}{n} \sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j | y_i) \quad 1 \leq j \leq m\end{aligned}$$

where

$$\pi_{\theta}(j | y_i) = \frac{\pi(j) \phi(y_i; \mu_j, \sigma_j^2)}{\sum_{j=1}^m \pi(j) \phi(y_i; \mu_j, \sigma_j^2)}, \quad 1 \leq j \leq m$$

E-step:

$$\begin{aligned}Q(\theta | \hat{\theta}^{(k)}) &= cst - \sum_{i=1}^n \sum_{j=1}^m \pi_{\hat{\theta}^{(k)}}(j | y_i) \frac{(y_i - \mu_j)^2}{2\sigma_j^2} \\ &= cst - \sum_{i=1}^n \sum_{j=1}^m \frac{\hat{\pi}^{(k)}(j) \phi(y_i; \hat{\mu}_j^{(k)}, \hat{\sigma}_j^{2(k)})}{\sum_{j=1}^m \hat{\pi}^{(k)}(j) \phi(y_i; \hat{\mu}_j^{(k)}, \hat{\sigma}_j^{2(k)})} \frac{(y_i - \mu_j)^2}{2\sigma_j^2}\end{aligned}$$

M-Step:

$$\begin{aligned}\hat{\mu}_j^{(k+1)} &= \frac{\sum_{i=1}^n y_i \pi_{\hat{\theta}^{(k)}}(j | y_i)}{\sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j | y_i)} \\ (\hat{\sigma}_j^2)^{(k+1)} &= \frac{\sum_{i=1}^n (y_i - \hat{\mu}_j^{(k+1)})^2 \pi_{\hat{\theta}^{(k)}}(j | y_i)}{\sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j | y_i)} \\ \hat{\pi}^{(k+1)}(j) &= \frac{1}{n} \sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j | y_i), \quad 1 \leq j \leq m.\end{aligned}$$

10.4 Convergence of EM Algorithm

Theorem 10.1

The likelihood sequence $p_{\hat{\theta}^{(k)}}(y)$, $k = 0, 1, 2, \dots$ is nondecreasing.



Proof 10.1

Assume for notational simplicity that the random variables Y and Z are discrete. Hence, their joint

distribution is given by

$$\begin{aligned} P_\theta(y, z) &= q_\theta(z)p_\theta(y|z) = q_\theta(z)\mathbf{1}_{\{y=h(z)\}} \\ &= p_\theta(y)q_\theta(z|y) \end{aligned}$$

Given y , the following identity holds for all $z \in h^{-1}(y)$:

$$p_\theta(y) = \frac{q_\theta(z)}{q_\theta(z|y)}$$

Taking the logarithm,

$$\ln p_\theta(y) = \ln q_\theta(z) - \ln q_\theta(z|y), \quad \forall z \in h^{-1}(y)$$

Taking the conditional expectation with respect to $q_{\hat{\theta}}(z|y)$,

$$\ln p_\theta(y) = \sum_{z \in h^{-1}(y)} q_{\hat{\theta}}(z|y) \ln q_\theta(z) - \sum_{z \in h^{-1}(y)} q_{\hat{\theta}}(z|y) \ln q_\theta(z|y) \quad (1)$$

Expectation (E)-Step: Compute

$$Q(\theta|\hat{\theta}^{(k)}) = \sum_{z \in h^{-1}(y)} q_{\hat{\theta}^{(k)}}(z|y) \ln q_\theta(z)$$

Maximization (M)-Step

$$\hat{\theta}^{(k+1)} = \underset{\theta \in S}{\operatorname{argmax}} Q(\theta|\hat{\theta}^{(k)})$$

According to (1),

$$\ln p_\theta(y) = Q(\theta|\hat{\theta}^{(k)}) - H(q_{\hat{\theta}^{(k)}}, q_\theta)$$

$$\ln p_{\hat{\theta}^{(k+1)}}(y) - \ln p_{\hat{\theta}^{(k)}}(y) = (Q(\hat{\theta}^{(k+1)}|\hat{\theta}^{(k)}) - Q(\hat{\theta}^{(k)}|\hat{\theta}^{(k)})) - (H(q_{\hat{\theta}^{(k)}}, q_{\hat{\theta}^{(k+1)}}) - H(q_{\hat{\theta}^{(k)}}, q_{\hat{\theta}^{(k)}}))$$

Since $\hat{\theta}^{(k+1)} = \underset{\theta \in S}{\operatorname{argmax}} Q(\theta|\hat{\theta}^{(k)})$, $Q(\hat{\theta}^{(k+1)}|\hat{\theta}^{(k)}) - Q(\hat{\theta}^{(k)}|\hat{\theta}^{(k)}) \geq 0$.

$$H(q_{\hat{\theta}^{(k)}}, q_{\hat{\theta}^{(k+1)}}) - H(q_{\hat{\theta}^{(k)}}, q_{\hat{\theta}^{(k)}}) = D(q_{\hat{\theta}^{(k)}} \| q_{\hat{\theta}^{(k+1)}}) \geq 0$$

Hence, we can conclude $\ln p_{\hat{\theta}^{(k+1)}}(y) - \ln p_{\hat{\theta}^{(k)}}(y) \geq 0$. Then $p_{\hat{\theta}^{(k)}}$ should be nondecreasing in k .

Corollary 10.1

Assume that S is a closed, bounded subset of Euclidean space, the functions $Q(\theta|\theta')$ and $H(\theta|\theta')$ are continuously differentiable, and the loglikelihood function $\ln p_{\hat{\theta}^{(k)}}$ is differentiable and bounded. Then the sequence $\ln p_{\hat{\theta}^{(k)}}$ converges, and any limit point $\theta^* \in \text{interior}(S)$ of the EM sequence is a solution of the likelihood equation $\nabla \ln p_\theta = 0$.



10.5 EM As an Alternating Maximization Algorithm

Define an auxiliary cost function $L(q, \theta)$.

Incomplete data Y ; Complete data Z . Still $h(z) : Z \rightarrow Y$.

$$\mathcal{Q}_y = \{q : q(z) = 0, \forall z \in h^{-1}(y)\}$$

EM updates

1. E-Step:

Chapter 11 Hidden Markov model (HMM)

A Markov chain $X_{t \geq 1}$ is observed as $\{Y_t\}_{t \geq 1}$. The state sets are finite sets S_x, S_y . Suppose the initial state distribution is π . The *transition probability matrix* of the MC is

$$A(i, j) = P(X_{t+1} = j | X_t = i), \quad i, j \in S_x$$

and the *emission probability matrix* is

$$B(i, j) = P(Y_t = j | X_t = i), \quad i \in S_x, j \in S_y$$

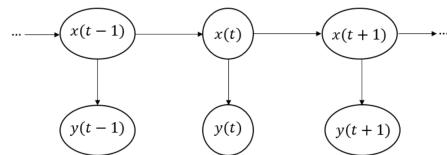


Figure 11.1: Hidden Markov Model (HMM)

Relative problems include

Problem 1: Estimate X_t given $Y_{1:t}$ (Using MAP or MMSE criterion: particle filtering)

Problem 2: Estimate X_{t+1} given $Y_{1:t}$ (Using MAP or MMSE prediction: particle filtering)

Problem 3: Estimate $X_{1:t}$ given $Y_{1:t}$ (MAP, MMSE)

Problem 4: Estimate the HMM parameters $\theta = (\pi, A, B)$ given $Y_{1:t}$ (learning)

11.1 Viterbi Algorithm: (MAP) estimate $X_{1:t}$ given $Y_{1:t}$

11.1.1 MAP estimation problem

The MAP estimation problem arises in a variety of applications, and Viterbi derived a remarkable algorithm for solving it exactly. The probability of state $\vec{x} \in S_x^n$ is given by

$$P(\vec{x}) = \pi(x_1) \prod_{t=1}^{n-1} A(x_t, x_{t+1})$$

and the conditional probability of the observed sequence \vec{y} given the state sequence \vec{x} is

$$P(\vec{y}|\vec{x}) = \prod_{t=1}^n B(x_t, y_t)$$

Hence, the joint probability of \vec{x} and \vec{y} is

$$P(\vec{x}, \vec{y}) = P(\vec{x})P(\vec{y}|\vec{x}) = \pi(x_1) \prod_{t=1}^{n-1} A(x_t, x_{t+1}) \prod_{t=1}^n B(x_t, y_t)$$

Then the MAP estimation problem is

$$\begin{aligned}\vec{x}^* &= \operatorname{argmax}_{\vec{x}} P(\vec{x}|\vec{y}) = \operatorname{argmax}_{\vec{x}} \frac{P(\vec{x}, \vec{y})}{P(\vec{y})} = \operatorname{argmax}_{\vec{x}} P(\vec{x}, \vec{y}) \\ &= \operatorname{argmax}_{\vec{x}} \pi(x_1) \prod_{t=1}^{n-1} A(x_t, x_{t+1}) \prod_{t=1}^n B(x_t, y_t) \\ &= \operatorname{argmax}_{\vec{x}} \ln \pi(x_1) + \sum_{t=1}^{n-1} \ln A(x_t, x_{t+1}) + \sum_{t=1}^n \ln B(x_t, y_t)\end{aligned}$$

11.1.2 Viterbi Algorithm

Let $f(x_1) = \ln \pi(x_1) + \ln B(x_1, y_1)$, $g_t(x_t, x_{t+1}) = \ln A(x_t, x_{t+1}) + \ln B(x_{t+1}, y_{t+1})$. Then the estimation problem is written in the form

$$\vec{x}^* = \operatorname{argmax}_{\vec{x}} \left[\varepsilon(\vec{x}) = f(x_1) + \sum_{u=1}^{n-1} g_u(x_u, x_{u+1}) \right]$$

Let $V(1, x) = f(x)$ and

$$V(t, x_t = x) \triangleq \max_{x_1, x_2, \dots, x_{t-1}} \left[\varepsilon([x_1, \dots, x_t]) = f(x_1) + \sum_{u=1}^{t-1} g_u(x_u, x_{u+1}) \right]$$

$$V(t, x_t = x) = \max_{x'} [V(t-1, x_{t-1} = x') + g_{t-1}(x', x)], t \geq 2$$

Then, when $t = n$ we have

$$\max_{\vec{x}} \varepsilon(\vec{x}) = \max_x V(n, x_n = x)$$

The complexity of the algorithm is $O(n|S_x|)$ storage and $O(n|S_x|^2)$ computation.

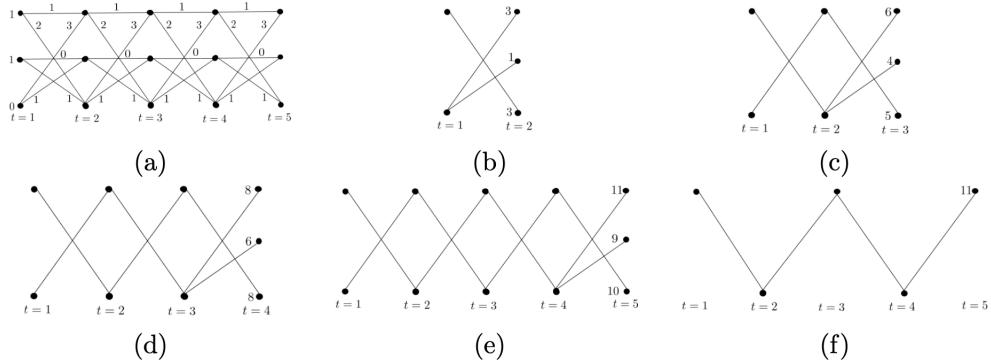


Figure 11.2: (a) Trellis diagram; (b)–(e) evolution of the Viterbi algorithm, showing surviving paths and values $V(t, x)$ at times $t = 2, 3, 4, 5$; (f) optimal path $\vec{x}^* = (0, 2, 0, 2, 0)$ and its value $\varepsilon(\vec{x}^*) = 11$.

11.2 Bayesian Estimation of a Sequence: Need (MMSE) estimate $X_{1:t}$ given $Y_{1:t}$

Consider Bayesian estimation under an additive squared-error loss function:

$$L(\vec{x}, \hat{\vec{x}}) = \sum_{t=1}^n L(x_t, \hat{x}_t) = \sum_{t=1}^n (x_t - \hat{x}_t)^2$$

The Bayesian estimator $\hat{\vec{x}}$ achieves

$$\min_{\hat{\vec{x}}} \sum_{\vec{x} \in X^n} L(\vec{x}, \hat{\vec{x}}) P(\vec{x} | \vec{y}) = \sum_{t=1}^n \min_{\hat{x}_t} \sum_{x_t \in X} L(x_t, \hat{x}_t) P(x_t | \vec{y})$$

In particular, under squared-error loss, we obtain the conditional mean estimator

$$\hat{x}_t = \sum_{x_t \in X} x_t P(X_t = x | Y = \vec{y}), \quad 1 \leq t \leq n$$

11.3 Forward-Backward Algorithm: (MMSE) estimate $X_{1:t+1}$ given $Y_{1:t}$

1. Evaluate $P(X_t = x | Y = \vec{y})$ for $t = 1, 2, \dots, n$ and $x \in \mathcal{X}$. (Used to (MMSE) estimate $X_{1:t}$ given $Y_{1:t}$)
2. Evaluate $P(X_t = x, X_{t+1} = x' | Y = \vec{y})$ for $t = 1, 2, \dots, n$ and $x, x' \in \mathcal{X}$ (Used to learn parameters $\theta = (\pi, A, B)$)

Define the shorthands

$$\begin{aligned} \gamma_t(x) &\triangleq P\left\{X_t = x \mid \vec{Y} = \vec{y}\right\} \\ \xi_t(x, x') &\triangleq P\left\{X_t = x, X_{t+1} = x' \mid \vec{Y} = \vec{y}\right\}, \quad x, x' \in \mathcal{X} \end{aligned}$$

Hence γ_t is the first marginal of ξ_t . The forward-backward algorithm allows efficient computation of these probabilities.

$$11.3.1 \quad \gamma_t(x) \triangleq P\left\{X_t = x \mid \vec{Y} = \vec{y}\right\}$$

We begin with

$$\gamma_t(x) = P\left\{X_t = x \mid \vec{Y} = \vec{y}\right\} = \frac{P\left\{X_t = x, \vec{Y} = \vec{y}\right\}}{\sum_{x \in \mathcal{X}} P\left\{X_t = x, \vec{Y} = \vec{y}\right\}}, \quad 1 \leq t \leq n$$

Write the numerator as a product of two conditional distributions,

$$\begin{aligned} P\left\{\vec{Y} = \vec{y}, X_t = x\right\} &\stackrel{(a)}{=} \underbrace{P\left\{Y_{1:t} = y_{1:t}, X_t = x\right\}}_{\mu_t(x)} \underbrace{P\left\{Y_{t+1:n} = y_{t+1:n} \mid X_t = x\right\}}_{\nu_t(x)} \\ &= \mu_t(x)\nu_t(x), \quad 1 \leq t < n \end{aligned}$$

where (a) follows from the Markov chain $Y_{1:t} \rightarrow X_t \rightarrow Y_{t+1:n}$. For $t = n$, we let $\nu_n(x) \equiv 1$. Combining above two equations we have

$$\gamma_t(x) = \frac{\mu_t(x)\nu_t(x)}{\sum_{x \in \mathcal{X}} \mu_t(x)\nu_t(x)}.$$

- (1) The first factor in the product of $P\left\{\vec{Y} = \vec{y}, X_t = x\right\}$ is

$$\mu_t(x) = P\{Y_{1:t} = y_{1:t}, X_t = x\}, \quad x \in \mathcal{X}, 1 \leq t \leq n,$$

for which we derive a **forward recursion**. The recursion is initialized with

$$\mu_1(x) = P\{Y_1 = y_1, X_1 = x\} = \pi(x)B(x, y_1).$$

For $t \geq 1$ we express μ_{t+1} in terms of μ_t as follows:

$$\begin{aligned} \mu_{t+1}(x) &= P\{Y_{1:t+1} = y_{1:t+1}, X_{t+1} = x\} \\ &\stackrel{(a)}{=} P\{Y_{1:t} = y_{1:t}, X_{t+1} = x\} P\{Y_{t+1} = y_{t+1} \mid X_{t+1} = x\} \\ &= B(x, y_{t+1}) \sum_{x' \in \mathcal{X}} P\{Y_{1:t} = y_{1:t}, X_{t+1} = x, X_t = x'\} \\ &\stackrel{(b)}{=} B(x, y_{t+1}) \sum_{x' \in \mathcal{X}} P\{Y_{1:t} = y_{1:t}, X_t = x'\} P\{X_{t+1} = x \mid X_t = x'\} \\ &= B(x, y_{t+1}) \sum_{x' \in \mathcal{X}} \mu_t(x') A(x', x), \quad t = 1, 2, \dots, n-1 \end{aligned}$$

where (a) holds because $Y_{1:t} \rightarrow X_{t+1} \rightarrow Y_{t+1}$ forms a Markov chain, and (b) because $Y_{1:t} \rightarrow X_t \rightarrow X_{t+1}$ forms a Markov chain.

- (2) The second factor in the product of $P\left\{\vec{Y} = \vec{y}, X_t = x\right\}$ is

$$\nu_t(x) = P\{Y_{t+1:n} = y_{t+1:n} \mid X_t = x\}, \quad x \in \mathcal{X}, 1 \leq t < n.$$

Starting from $\nu_n(x) \equiv 1$, we have the following **backward recursion**, expressing ν_{t-1} in terms of ν_t for $2 \leq t \leq n$:

$$\begin{aligned} \nu_{t-1}(x) &= P\{Y_{t:n} = y_{t:n} \mid X_{t-1} = x\} \\ &= \sum_{x' \in \mathcal{X}} P\{Y_{t:n} = y_{t:n}, X_t = x' \mid X_{t-1} = x\} \\ &\stackrel{(a)}{=} \sum_{x' \in \mathcal{X}} P\{Y_{t:n} = y_{t:n} \mid X_t = x'\} P\{X_t = x' \mid X_{t-1} = x\} \\ &\stackrel{(b)}{=} \sum_{x' \in \mathcal{X}} P\{Y_{t+1:n} = y_{t+1:n} \mid X_t = x'\} P\{Y_t = y_t \mid X_t = x'\} P\{X_t = x' \mid X_{t-1} = x\} \\ &= \sum_{x' \in \mathcal{X}} \nu_t(x') B(x', y_t) A(x, x'), \quad t = n, n-1, \dots, 2 \end{aligned}$$

where (a) holds because $X_{t-1} \rightarrow X_t \rightarrow Y_{t:n}$ forms a Markov chain, and (b) because $Y_{t+1:n} \rightarrow X_t \rightarrow Y_t$ forms a Markov chain.

$$11.3.2 \quad \xi_t(x, x') \triangleq P\left\{X_t = x, X_{t+1} = x' \mid \vec{Y} = \vec{y}\right\}$$

Next we derive an expression for

$$\xi_t(x, x') = P\left\{X_t = x, X_{t+1} = x' \mid \vec{Y} = \vec{y}\right\} = \frac{P\left\{\vec{Y} = \vec{y}, X_t = x, X_{t+1} = x'\right\}}{\sum_{x, x' \in \mathcal{X}} P\left\{\vec{Y} = \vec{y}, X_t = x, X_{t+1} = x'\right\}}$$

We have

$$\begin{aligned} & P\left\{\vec{Y} = \vec{y}, X_t = x, X_{t+1} = x'\right\} \\ & \stackrel{(a)}{=} P\left\{Y_{1:t+1} = y_{1:t+1}, X_t = x, X_{t+1} = x'\right\} P\left\{Y_{t+2:n} = y_{t+2:n} \mid X_{t+1} = x'\right\} \\ & \stackrel{(b)}{=} P\left\{Y_{1:t} = y_{1:t}, X_t = x\right\} P\left\{X_{t+1} = x' \mid X_t = x\right\} P\left\{Y_{t+1} = y_{t+1} \mid X_{t+1} = x'\right\} \nu_{t+1}(x') \\ & = \mu_t(x) A(x, x') B(x', y_{t+1}) \nu_{t+1}(x') \end{aligned}$$

where (a) holds because $(Y_{1:t+1}, X_t) \rightarrow X_{t+1} \rightarrow Y_{t+2:n}$ forms a Markov chain, and (b) because $Y_{1:t} \rightarrow X_t \rightarrow X_{t+1} \rightarrow Y_{t+1}$ forms a Markov chain. Hence

$$\xi_t(x, x') = \frac{\mu_t(x) A(x, x') B(x', y_{t+1}) \nu_{t+1}(x')}{\sum_{x, x' \in \mathcal{X}} \mu_t(x) A(x, x') B(x', y_{t+1}) \nu_{t+1}(x')}, \quad 1 \leq t \leq n, x, x' \in \mathcal{X}$$

11.3.3 Scaling Factors

Unfortunately the recursions above are numerically unstable for large n because the probabilities $\mu_t(x)$ and $\nu_t(x)$ vanish exponentially with n and are sums of many small terms of different sizes. The following approach is more stable. Define

$$\begin{aligned} \alpha_t(x) &= P\left\{X_t = x \mid Y_{1:t} = y_{1:t}\right\}, \\ \beta_t(x) &= \frac{P\left\{Y_{t+1:n} = y_{t+1:n} \mid X_t = x\right\}}{P\left\{Y_{t+1:n} = y_{t+1:n} \mid Y_{1:t} = y_{1:t}\right\}}, \\ c_t &= P\left\{Y_t = y_t \mid Y_{1:t-1} = y_{1:t-1}\right\} \end{aligned}$$

Then

$$\gamma_t(x) = \alpha_t(x) \beta_t(x)$$

$$\xi_t(x, x') = c_t \alpha_t(x) B(x, y_t) A(x, x') \beta_t(x')$$

A forward recursion can be derived for α_t and c_t , and a backward recursion for β_t .

The time and storage complexity of the algorithm is $O(n|\mathcal{X}|^2)$.

Chapter 12 Graphic Models

To compute $P(x_1, \dots, x_d)$, we can utilize the chain rule $P(x_1, \dots, x_d) = P(x_1) \prod_{i=2}^d P(x_i | x_{1:i-1})$. However, this approach becomes computationally expensive as the dimension d increases.

Fortunately, when there are conditionally independent relationships between variables, such as $x_A \perp x_C | x_B$, we can reduce the computational cost.

In this section, we can employ graphical models to represent probabilistic relationships between variables, particularly when there are conditionally independent relationships present.

12.1 Graph Theory

1. A graph (V, E) , V is a set of *vertices*, $E \subseteq V \times V$ is a set of ordered pairs of vertices, called *edges*.

An edge $(i, j) \in E$ is *directed* if $(i, j) \notin E$; otherwise the edge is *undirected*. We denote directed and undirected edges by the symbols $i \rightarrow j$ and $i \sim j$, respectively.

2. **Directed and Undirected Graphs:** Graphs in which *all* edges are directed (resp. undirected).
3. **Subgraph:** a subgraph (S, E_S) of (V, E) is a subset $S \subset G$ with edges that have both endpoints in S .
4. **Clique:** A set C of vertices in an undirected graph is a clique if either C is a singleton, or **each pair of vertices in C is linked by an edge**.

That is, all vertices in C are neighbors. The clique is maximal if there is no larger clique that contains C .

5. **Parent, Child:** Vertex i is a parent of vertex j if $i \rightarrow j$, in which case j is also called a child of i . We denote by $\pi(j)$ the set of parents of j .

6. **Path:** A *path* of length n from i to j is a sequence $i = k_0, k_1, \dots, k_n = j$ of distinct vertices such that $(k_{m-1}, k_m) \in E$ for all $m = 1, \dots, n$. We designate such a path by $i \rightarrow j$.

7. **Connected Graph:** An undirected graph is connected if there is a path between any pair of nodes. In general, the connected components of a graph are those subgraphs which are connected.

8. **Cycle/Loop:** An n -cycle, or loop, is a path of length n $i \rightarrow j$ with $i = j$.

A directed graph without cycle is also called Directed Acyclic Graph (DAG)

9. **Tree:** A tree is a connected, undirected graph without cycles; **it has a unique path between any two vertices**.

10. **Rooted Tree:** A rooted tree is the directed acyclic graph obtained from a tree by choosing as vertex as root and directing all edges away from this root. Each vertex of a rooted tree has at most one parent.

11. **Forest:** A forest is an undirected graph where all connected components are trees.

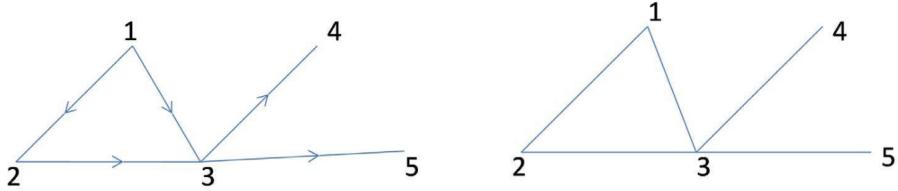


Figure 12.1: (a) Directed and (b) Undirected graph.

12.2 Bayesian Networks

A Bayesian network (or belief network) is a joint probability distribution associated with a *directed acyclic graph* (V, E) whose nodes $X_v, v \in V$ are random variables. The joint distribution is of the form

$$p(\vec{x}) = \prod_{v \in V} p(x_v | \pi(x_v))$$

$\pi(x_v)$ is the set of parents of vertices.

For instance a Markov chain is a chain-type directed acyclic graph where $V = \{1, 2, \dots, n\}$, and $\pi(v) = v - 1$ for $v \geq 2$. The pmf for the sequence \vec{x} is obtained from the chain rule

$$p(\vec{x}) = p(x_1)p(x_2|x_1) \cdots p(x_n|x_{n-1})$$

12.3 Markov Networks

12.3.1 General Form

We can use undirected graph to represent conditionally independent.

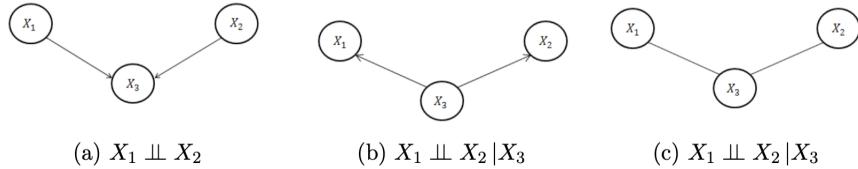


Figure 12.2: (a) (b) Two Bayesian networks and (c) a Markov network.

More generally, if two nodes X_u and X_v in a Markov network are not connected by an edge, then the random variables X_u and X_v are conditionally independent given all the other random variables (denoted by $X_u \perp X_v | X_{V \setminus \{u,v\}}$).

A Markov network is an undirected graph $G = (V, E)$ together with a collection $X = \{X_v, v \in V\}$ of random variables indexed by the nodes of G .

Since there is no direction, we use **clique** to help use represent probabilities. (*Review: clique is a set of vertices that each pair of vertices is linked*)

We use Ω let be collection of cliques in the graph and the functions $\psi_C(\cdot)$ be the **clique potentials, or compatibility functions**.

The pmf of X takes the form

$$p(\vec{x}) = \frac{\prod_{C \in \Omega} \psi_C(\vec{x}_C)}{\sum_{\vec{x}} \prod_C \psi_C(\vec{x}_C)} = \frac{1}{Z} \prod_{C \in \Omega} \psi_C(\vec{x}_C)$$

where $Z = \sum_{\vec{x}} \prod_C \psi_C(\vec{x}_C)$ is a normalization constant.

Note: this is a form of factorization that can represent conditionally independent relationship among variables.

$\psi_C(\cdot)$ are undefined functions.x

12.3.2 Hammersley-Clifford theorem

Theorem 12.1 (Hammersley-Clifford theorem)

Assume that $p(x_1, \dots, x_n) > 0$ (positivity condition). Then,

$$p(\vec{x}) = \frac{1}{Z} \prod_{C \in \Omega} \phi_C(\vec{x}_C)$$

Thus, the following are equivalent (given the positivity condition):

1. **Local Markov property:** $p(x_i | \vec{x} \setminus \{x_i\}) = p(x_i | \mathcal{N}(x_i))$, where $\mathcal{N}(x_i)$ is the neighboring set of x_i .
2. **Factorization property:** The probability factorizes according to the cliques of the graph.
3. **Global Markov property:** $p(\vec{x}_A | \vec{x}_B, \vec{x}_S) = p(\vec{x}_A | \vec{x}_S)$ whenever \vec{x}_A and \vec{x}_B are separated by \vec{x}_S in G



12.3.3 Form of Gibbs distribution (Boltzmann distribution)

The factorization is not unique. We let $\psi(\vec{x}_C) = e^{-V_C(\vec{x}_C)}$, where $V_C(\cdot)$ are the so-called potential energy functions. In a pairwise Markov network, $p(\vec{x})$ can be expressed as a product of clique potentials involving either one or two random variables.

$$p(\vec{x}) = \frac{1}{Z} e^{-\sum_C V_C(x_C)}$$

This probability follows **Gibbs distribution (Boltzmann distribution)**. This distribution follows exponential families.

12.4 Conversion of directed graph to undirected graph

We can use a step known as *moralization*. Moralization of graph: connect two unmarried parents.

This is the process of “marrying” the parents of each node, i.e., adding an edge connecting any pair of parents if one did not exist. The figure illustrates this process for a node with three parents. In this case the undirected graph consists of a clique of size 4.

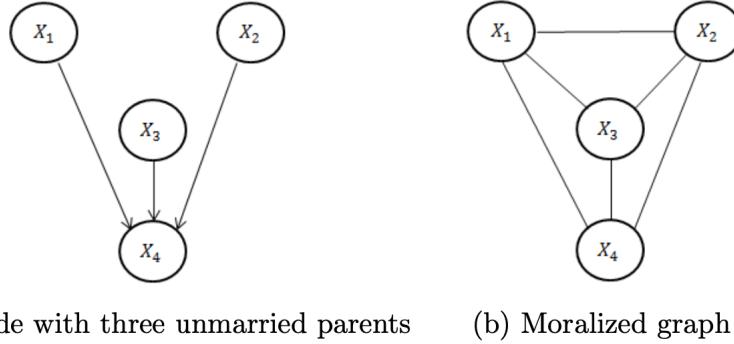


Figure 12.3: Graph moralization

12.5 Inference and Learning

12.5.1 Inference on Trees

Consider the tree of the figure, which has 5 nodes and edges $1 \sim 2 \sim 3$ and $4 \sim 3 \sim 5$. We have

$$p(\vec{x}) = \frac{1}{Z} \psi_{12}(x_1, x_2) \psi_{23}(x_2, x_3) \psi_{34}(x_3, x_4) \psi_{35}(x_3, x_5)$$

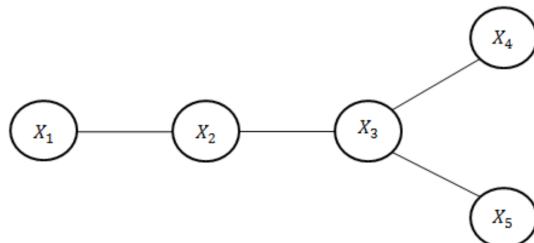


Figure 12.4: Example 1

1. **Marginal Inference:** As a first example of inference on trees, consider the problem of evaluating the marginal pmf $p(x_5)$. We explore two approaches: the **direct approach**, which is computationally

infeasible for large graphs (the number of items in the sum is $|\mathcal{X}|^4$);

$$p(x_5) = \sum_{x_1, \dots, x_4} \frac{1}{Z} \psi_{12}(x_1, x_2) \psi_{23}(x_2, x_3) \psi_{34}(x_3, x_4) \psi_{35}(x_3, x_5)$$

and the **sum-product algorithm**, which exploits the graph structure.

$$\begin{aligned} p(x_5) &= \frac{1}{Z} \sum_{x_3} \psi_{35}(x_3, x_5) \sum_{x_4} \psi_{34}(x_3, x_4) \sum_{x_2} \underbrace{\psi_{23}(x_2, x_3) \sum_{x_1} \psi_{12}(x_1, x_2)}_{m_{1 \rightarrow 2}(x_2)} \\ &= \frac{1}{Z} \sum_{x_3} \psi_{35}(x_3, x_5) \sum_{x_4} \psi_{34}(x_3, x_4) \underbrace{\sum_{x_2} \psi_{23}(x_2, x_3) m_{1 \rightarrow 2}(x_2)}_{m_{2 \rightarrow 3}(x_3)} \\ &= \frac{1}{Z} \sum_{x_3} \psi_{35}(x_3, x_5) m_{2 \rightarrow 3}(x_3) \underbrace{\sum_{x_4} \psi_{34}(x_3, x_4)}_{m_{4 \rightarrow 3}(x_3)} \\ &= \frac{1}{Z} \underbrace{\sum_{x_3} \psi_{35}(x_3, x_5) m_{2 \rightarrow 3}(x_3) m_{4 \rightarrow 3}(x_3)}_{m_{3 \rightarrow 5}(x_5)}. \end{aligned}$$

In this derivation, nodes 1, 2, 4, 3 are eliminated in that order. We think of each term $m_{i \rightarrow j}(x_j)$ as a message conveyed from node i to node j , just before elimination of j . Computing $m_{i \rightarrow j}(x_j)$ involves a summation over all possible values of x_i . This interpretation will be helpful in more complex problems.

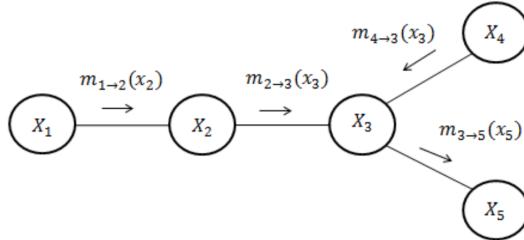


Figure 12.5: Belief propagation in a tree

As illustrated in the Figure, a node can send a message to a neighbor once it has received messages from all of its other neighbors. For a general tree, upon choosing an elimination order, we evaluate the following messages in the corresponding order:

$$m_{i \rightarrow j}(x_j) = \sum_{x_i} \psi_{ij}(x_i, x_j) \prod_{k \in \mathcal{N}(i) \setminus \{j\}} m_{k \rightarrow i}(x_i)$$

The marginal probability at any node i is the product of all incoming messages:

$$p(x_i) = \frac{1}{Z} \prod_{k \in \mathcal{N}(i)} m_{k \rightarrow i}(x_i).$$

We can also evaluate the 2D marginal $p(x_2, x_5)$

$$p(x_2, x_5) = \frac{1}{Z} \sum_{x_3} \underbrace{\psi_{23}(x_2, x_3)}_{m_{4 \rightarrow 3}(x_3)} \underbrace{\psi_{35}(x_3, x_5)}_{m_{1 \rightarrow 2}(x_2)} \underbrace{\sum_{x_4} \psi_{34}(x_3, x_4)}_{\overbrace{x_1}^{\psi_{12}(x_1, x_2)}} \sum_{x_1} \psi_{12}(x_1, x_2).$$

Finally, a conditional marginal such as $p(x_1 | x_5)$ is obtained as $p(x_1, x_5) / p(x_5)$, hence the problem is reduced to evaluating unconditional marginals.

The computational cost of the algorithm is $O(n|\mathcal{X}|^2)$ when the n random variables are defined over the same alphabet \mathcal{X} .

2. **Maximization:** A closely related problem is to find the most likely configuration, possibly by fixing some coordinates. For instance, evaluate

$$M(x_5) = \max_{x_1, x_2, x_3, x_4} p(x_1, x_2, x_3, x_4, x_5)$$

for the above Markov network. Direct calculation has exponential complexity. However, the more efficient max-product algorithm has the same structure as the sum-product algorithm:

$$\begin{aligned} M(x_5) &= \max_{x_1, x_2, x_3, x_4} \frac{1}{Z} \psi_{12}(x_1, x_2) \psi_{23}(x_2, x_3) \psi_{34}(x_3, x_4) \psi_{35}(x_3, x_5) \\ &= \frac{1}{Z} \max_{x_3} \psi_{35}(x_3, x_5) \max_{x_4} \psi_{34}(x_3, x_4) \max_{x_2} \psi_{23}(x_2, x_3) \underbrace{\max_{x_1} \psi_{12}(x_1, x_2)}_{m_{1 \rightarrow 2}(x_2)} \\ &= \frac{1}{Z} \max_{x_3} \psi_{35}(x_3, x_5) \underbrace{\max_{x_4} \psi_{34}(x_3, x_4)}_{m_{4 \rightarrow 3}(x_3)} \underbrace{\max_{x_2} \psi_{23}(x_2, x_3)}_{m_{2 \rightarrow 3}(x_3)} \underbrace{\max_{x_1} \psi_{12}(x_1, x_2)}_{m_{1 \rightarrow 2}(x_2)} \\ &= \frac{1}{Z} \underbrace{\max_{x_3} \psi_{35}(x_3, x_5)}_{m_{3 \rightarrow 5}(x_5)} \underbrace{\max_{x_4} \psi_{34}(x_3, x_4)}_{m_{4 \rightarrow 3}(x_3)} \underbrace{\max_{x_2} \psi_{23}(x_2, x_3)}_{m_{2 \rightarrow 3}(x_3)} \end{aligned}$$

Chapter 13 Variational Inference, Mean-Field Techniques

Approximate complicated p.d.f. $p(\vec{x})$ with tractable $q(\vec{x})$, where $q \in Q =$ tractable set of distributions.

Use divergence to measure:

$$\min_{q \in Q} D(q \| p)$$

note that D is convex in q .

13.1 Naive Mean-Field Methods

The *naive mean field method* approximates a distribution by a product distribution.

Assume the q has the form $q(\vec{x}) = \prod_{i=1}^n q_i(x_i)$. Assume $x_i \in X =$ finite set.

$$\begin{aligned} D(q \| p) &= \mathbb{E}_q \left[\ln \frac{q(\vec{X})}{p(\vec{X})} \right] \\ &= \sum_{i=1}^n \mathbb{E}_q [\ln q_i(X_i)] - \mathbb{E}_q [\ln p(\vec{X})] \\ &= \sum_{i=1}^n \sum_{x_i \in X} q_i(x_i) \ln q_i(x_i) - \sum_{\vec{x} \in X^n} \left(\prod_{i=1}^n q_i(x_i) \right) \ln p(\vec{x}) \end{aligned}$$

Solve

$$\begin{aligned} \min_{\{q_i\}} \quad & D\left(\prod_{i=1}^n q_i \| p\right) \\ \text{s.t. } & \sum_{x_i \in X} q_i(x_i) = 1, i = 1, \dots, n \end{aligned}$$

Using Lagrangian method:

$$\begin{aligned} L(q, \vec{\lambda}) &= D(q \| p) + \sum_{i=1}^n \lambda_i \left(\sum_{x \in X} q_i(x_i) - 1 \right) \\ 0 &= \frac{\partial L(q, \vec{\lambda})}{\partial q_i(x_i)} = 1 + \ln q_i(x_i) - \sum_{\vec{x}' : x'_i = x_i} \left(\prod_{j \neq i}^n q_j(x'_j) \right) \ln p(\vec{x}') + \lambda_i \end{aligned}$$

Hence, $q_i(x_i)$ should in the form:

$$\begin{aligned} q_i(x_i) &= \frac{1}{e^{1+\lambda_i}} e^{\sum_{\vec{x}' : x'_i = x_i} \left(\prod_{j \neq i}^n q_j(x'_j) \right) \ln p(\vec{x}')} \\ &= \frac{1}{Z_i} \exp \left(\mathbb{E}_{\prod_{j \neq i}^n q_j} [\ln p(X_{1:i-1}, x_i, X_{i+1:n})] \right) \end{aligned}$$

Iteration Algorithm:

$$q_i^{(k+1)}(x_i) = \frac{1}{Z_i} \exp \left(\mathbb{E}_{\prod_{j \neq i}^n q_j^{(k)}} [\ln p(X_{1:i-1}, x_i, X_{i+1:n})] \right)$$

13.1.1 Graphical Models

Consider $P = \text{pairwise Markov model}$

$$p(\vec{x}) = \frac{1}{Z} \prod_{(i,j) \in E} \psi_{ij}(x_i, x_j)$$

$$\ln p(\vec{x}) = -\ln Z + \sum_{(i,j) \in E} \ln \psi_{ij}(x_i, x_j)$$

The expectation can be written as

$$\begin{aligned} \mathbb{E}_{\prod_{j \neq i} q_j} [\ln p(X_{1:i-1}, x_i, X_{i+1:n})] &= \mathbb{E}_{\prod_{j \neq i} q_j} \left[\sum_{(i,j) \in E} \ln \psi_{ij}(x_i, x_j) \right] + cst \\ &= \sum_{j \in N(i)} \mathbb{E}_{q_j} [\ln \psi_{ij}(x_i, x_j)] + cst \\ &= \sum_{j \in N(i)} \sum_{x_j \in X} q_j(x_j) \ln \psi_{ij}(x_i, x_j) + cst \end{aligned}$$

and thus,

$$q_i(x_i) = \frac{1}{Z_i} \exp \left(\sum_{j \in N(i)} \sum_{x_j \in X} q_j(x_j) \ln \psi_{ij}(x_i, x_j) \right)$$

13.1.2 Ising Model

Consider a 2-D torus V with $|V| = n$ nodes, and $X = \{\pm 1\}$. Each node is connected to its upper, lower, right, and left neighbors. The distribution is of the form

$$p(\vec{x}) = \frac{1}{Z} \exp \left(\beta \sum_{(i,j) \in E} x_i x_j \right), \quad \vec{x} \in \{\pm 1\}^n$$

with $\beta \geq 0$. The parameter β represents the inverse of a temperature. For $\beta = 0$ the distribution is uniform, hence fully factorized. For large positive values of β , configurations \vec{x} with strong correlations are favored.

2-D Ising Model

$$\psi_{ij}(x_i, x_j) = e^{\beta x_i x_j}$$

$$\Rightarrow q_i(x_i) = \frac{1}{Z_i} \prod_{j \in N(i)} \exp \left(- \sum_{x_j = \pm 1} q_j(x_j) \beta x_i x_j \right)$$

Since each X_i is a Bernoulli random variable, the decision variable q_i can be represented by a single parameter which we choose to be the mean $m_i = q_i(1) - q_i(-1) \in [-1, 1]$. Equivalently,

$$\begin{aligned} m_i = q_i(1) - q_i(-1) \Leftrightarrow q_i(1) &= \frac{1 + m_i}{2}, \quad q_i(-1) = \frac{1 - m_i}{2} \\ \Rightarrow q_i(x_i) &= \frac{1}{Z_i} \prod_{j \in N(i)} \exp (-\beta x_i m_j) \end{aligned}$$

Then, our problem is finding the optimal $\{m_i\}$.

$$q_i(1) = \frac{1}{Z_i} \exp\left(-\beta \sum_{j \in N(i)} m_j\right); q_i(-1) = \frac{1}{Z_i} \exp\left(\beta \sum_{j \in N(i)} m_j\right)$$

The normalization constant is given by

$$Z_i = \exp\left(\beta \sum_{j \in N(i)} m_j\right) + \exp\left(-\beta \sum_{j \in N(i)} m_j\right) = 2 \cosh\left(\beta \sum_{j \in N(i)} m_j\right)$$

Hence,

$$m_i = q_i(1) - q_i(-1) = \tanh\left(\beta \sum_{j \in N(i)} m_j\right)$$

Convergence. We show that the algorithm always converges if a uniform initialization is used, i.e., $m_i^{(0)} = m^{(0)}$ for all $i \in V$. Then the (simultaneous) update equation for the means is

$$m^{(k+1)} = \tanh\left(4\beta m^{(k)}\right)$$

where the factor of 4 arises because each vertex has 4 neighbors. Analysis of convergence depends on the value of β , and we need consider two cases.

- (1) **Case I:** $\beta < \frac{1}{4}$: The mapping is a contraction mapping for $\beta < \frac{1}{4}$, and so the fixed point of this mapping is $\lim_{k \rightarrow \infty} m^{(k)} = 0$, for any initialization $m^{(0)}$. Hence, the variational approximation is uniform: $q(\vec{x}) = 2^{-n}$ for all $x \in \{\pm 1\}^V$.
- (2) **Case II:** $\beta > \frac{1}{4}$: In this case, the equation $m = \tanh(4\beta m)$ has three possible solutions 0 and $\pm m^*$ where $m^* > 0$. If the algorithm is initialized with $m^{(0)} = 0$, then subsequent iterations do not change this value. If the algorithm is initialized with $m^{(0)} > 0$, it converges to m^* . Finally, if the algorithm is initialized with $m^{(0)} < 0$, it converges to $-m^*$. In the latter two cases (convergence to either m^* or $-m^*$), the variational approximations q_i are nonuniform.
- (3) **Case III:** $\beta = \frac{1}{4}$: phase transition

The case $\beta > \frac{1}{4}$ is related to percolation theory in statistical physics. It may be shown that the distribution p favors configurations featuring large homogeneous regions. The correlation between any two nodes is significant, even for large graphs. This behavior is completely different from the case $\beta < \frac{1}{4}$, where the correlation between distant nodes dies out with distance (similarly to a homogeneous, irreducible Markov chain). The case $\beta = \frac{1}{4}$ is known as a phase transition.

13.2 Exponential Families of Probability Distributions

Definition 13.1 (d -dimensional exponential families)

In canonical form: d -dimensional exponential families

$$p_\theta(y) = \frac{h(y)}{Z(\theta)} e^{\sum_{k=1}^d \theta_k T_k(y)}$$

$T_k(\cdot)$ are called sufficient statistics. **Partition function** $Z(\theta)$ is the normalization constant ensuring that the density p_θ integrates to 1.

It can also be written

$$p_\theta(y) = e^{\theta^T T(y) - A(\theta)}$$

The **log partition function** (aka cumulant function) $A(\theta) = \ln Z(\theta)$



Example 1: $P_\theta = N(\theta, 1)$,

$$p_\theta(y) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-\theta)^2}{2}} = \frac{e^{\frac{y^2}{2}}}{\sqrt{2\pi} e^{-\frac{\theta^2}{2}}} e^{-\theta y}$$

Example 2: $P_\theta = N(0, \theta^{-1})$, θ is the inverse covariance matrix.

$$p_\theta(y) = \frac{|\theta|^{\frac{1}{2}}}{\sqrt{2\pi}} e^{-\frac{1}{2} y^T \theta y}$$

Example 3: 2D-Ising Model

$$p_\theta(y) = \frac{1}{Z(\theta)} e^{\theta y_i y_j}, \theta > 0$$

Generalized 2D-Ising Model

$$p_\theta(y) = \frac{1}{Z(\theta)} e^{\sum_{i \in V} \theta_i y_i + \sum_{i \sim j} \theta_{ij} y_i y_j}$$

The natural parameter set:

$$\Theta = \left\{ \theta : \int_X e^{\theta^T T(y)} dy < \infty \right\}$$

The divergence in exponential form

$$\begin{aligned} D(P_\theta \| P_{\theta'}) &= \mathbb{E}_\theta \left[\ln \frac{P_\theta(Y)}{P_{\theta'}(Y)} \right] \\ &= \mathbb{E}_\theta [\theta^T T(Y) - A(\theta) - (\theta')^T T(Y) + A(\theta')] \\ &= -[A(\theta) - A(\theta')] + (\theta - \theta')^T \mathbb{E}_\theta [T(Y)] \end{aligned}$$

Cumulant-generating function (cgf)

$$\kappa(u) = \ln \mathbb{E}(e^{u^T T(Y)}) = \ln \int e^{u^T T(y)} (e^{\theta^T T(y) - A(\theta)}) dy$$

$$\nabla \kappa(0) = \mathbb{E}_\theta[T(Y)]$$

$$\nabla^2 \kappa(0) = \text{Cov}_\theta[T(Y)]$$

We can compute

$$\begin{aligned} \int e^{\theta^T T(y) - A(\theta)} dy &= 1 \\ \int [T(y) - \nabla A(\theta)] e^{\theta^T T(y) - A(\theta)} dy &= 0 \\ \int [T(y) - \nabla A(\theta)] p_\theta(y) dy &= 0 \\ \mathbb{E}_\theta[T(Y)] &= \nabla A(\theta) \end{aligned}$$

Hence,

$$\nabla A(\theta) = \nabla \kappa(0) = \mathbb{E}_\theta[T(X)]$$

$$\nabla^2 A(\theta) = \nabla^2 \kappa(0) = \text{Cov}_\theta[T(X)]$$

Definition. The set of realizable mean parameters \mathcal{M} is the set of μ that are the expected value of $T(X)$ under some distribution p (not necessarily in the exponential family). Thus

$$\mathcal{M} \triangleq \left\{ \mu \in \mathbb{R}^d : \exists p : \mathbb{E}_p[T(X)] = \mu \right\}$$

which is a convex set.

Example 6. Consider Generalized 2D-Ising Model

$$p_\theta(y) = \frac{1}{Z(\theta)} e^{\sum_{i \in V} \theta_i y_i + \sum_{i \sim j} \theta_{ij} y_i y_j}$$

For $y \in \{0, 1\}$

$$\mu_i = \mathbb{E}_\theta[T_i(Y)] = P_\theta(Y_i = 1)$$

$$\mu_{ij} = \mathbb{E}_\theta[T_i(Y)T_j(Y)] = P_\theta(Y_i = Y_j = 1)$$

Example 7. If $T(x) = xx^\top \in \mathbb{R}^{n \times n}$ then μ is a correlation matrix, and so \mathcal{M} is the set of all $n \times n$ symmetric nonnegative definite matrices.

13.3 ML Estimation

Consider n iid samples $X^{(i)}, 1 \leq i \leq n$ drawn from the exponential distribution p_θ . The ML estimator of θ given these n samples is obtained by solving

$$\begin{aligned}\hat{\theta}_{ML} &= \max_{\theta} \frac{1}{n} \sum_{i=1}^n \ln p_\theta(X^{(i)}) \\ &= \max_{\theta} \frac{1}{n} \sum_{i=1}^n [\theta^\top T(X^{(i)}) - A(\theta)] \\ &= \max_{\theta} [\theta^\top \hat{\mu} - A(\theta)]\end{aligned}$$

where $\hat{\mu}$ is the mean parameter

$$\hat{\mu} \triangleq \frac{1}{n} \sum_{i=1}^n T(X^{(i)})$$

$A(\theta)$ is a convex function, we can solve optimal solution by solving critical point.

$$\nabla A(\hat{\theta}_{ML}) = \hat{\mu}$$

The gradient mapping could be hard to invert, however. For instance, for the Ising model example we easily obtain

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \sum_{j \sim k} X_j^{(i)} X_k^{(i)}$$

but inverting the gradient mapping is a hard problem. Such is generally the case if p is a distribution over a Markov network with cycles.

13.4 Maximum Entropy

Consider a random variable X over a finite set \mathcal{X} . Its probability distribution p is unknown, however we are given the expected value $\mu_k = \mathbb{E}_p [T_k(X)]$ of d statistics $T_k(X), 1 \leq k \leq d$. A classical problem, which originates from statistical physics, is to find p that maximizes the entropy $H(p) = -\sum_x p(x) \ln p(x)$ subject to the d constraints above. Assuming the feasible set is nonvoid, the resulting distribution is called the maximum-entropy (or maxent()) distribution.

Since $H(p)$ is concave, the constraints are linear in p , and the probability simplex is a convex set, the maxent problem is concave. Its solution is obtained by introducing d Lagrange multipliers $\lambda_k, 1 \leq k \leq d$ associated with the mean constraints, and a Lagrange multiplier λ_{d+1} associated with the constraint $\sum_x p(x) = 1$. Ignoring momentarily the nonnegativity constraints, we maximize the Lagrangian

$$\mathcal{L}(p, \lambda) \triangleq -\sum_{x \in \mathcal{X}} p(x) \ln p(x) + \sum_{k=1}^d \lambda_k \left(\sum_{x \in \mathcal{X}} p(x) T_k(x) - \mu_k \right) + \lambda_{d+1} \left(\sum_{x \in \mathcal{X}} p(x) - 1 \right)$$

over p , subject to the $d + 1$ equality constraints

The first-order optimality conditions are given by

$$0 = \frac{\partial \mathcal{L}(p, \lambda)}{\partial p(x)} = -\ln p(x) - 1 + \sum_{k=1}^d \lambda_k T_k(x) + \lambda_{d+1}$$

Hence,

$$p(x) = \frac{1}{Z} e^{\sum_{k=1}^d \lambda_k T_k(x)}$$

where $Z = e^{1-\lambda_{d+1}}$

$$H(p) = \mathbb{E}_p[\ln p(X)] = -\theta^T \mathbb{E}_p[T(X)] + A(\theta) = -\max_{\theta} (\theta^T \mu - A(\theta))$$

Example. Let $X = (X_1, X_2) \in \{0, 1\}^2$ and consider maximizing entropy subject to the constraint $\mathbb{E}[X_1 X_2] = \mu$ where $\mu \in (0, 1)$. We obtain $p(x) = \frac{1}{Z} \exp\{\lambda x_1 x_2\}$. Since $\sum_x p(x) = 1$, the normalization constant is obtained as $Z = e^\lambda + 3$. We obtain λ from the constraint

$$\mu = \mathbb{E}_P[X_1 X_2] = \frac{e^\lambda}{e^\lambda + 3} \Rightarrow \lambda = \ln \frac{3\mu}{1-\mu}$$

The maxent solution takes the form

$$p(x) = \begin{cases} \mu & (x_1, x_2) = (1, 1) \\ \frac{1-\mu}{3} & \text{else} \end{cases}$$

and has entropy is $H(p) = -\mu \ln \mu - (1-\mu) \ln \frac{1-\mu}{3}$.

A similar version of the maxent problem exists for continuous random variables. The entropy function is replaced with the differential entropy functional $h(p) \triangleq -\int p \ln p$, and the maxent solution again takes an exponential form.

13.5

$$\begin{aligned} \min_{q \in Q} D(q||p) &= \min_q \mathbb{E}_q \left[\ln \frac{q(x)}{p_\theta(x)} \right] \\ &= \min_q [A(\theta) - \theta^T \mathbb{E}_q[T(x)] - H(q)] \\ &= A(\theta) - \max_{\mu} \max_{q: \mathbb{E}_q[T(x)] = \mu} [\theta^T \mathbb{E}_q[T(x)] + H(q)] \\ &= A(\theta) - \max_{\mu \in M} \left[\theta^T \mu + \max_{q: \mathbb{E}_q[T(x)] = \mu} H(q) \right] \end{aligned}$$

Since $\max_{q: \mathbb{E}_q[T(x)] = \mu} H(q)$ is exactly an entropy maximum problem, we let $A^*(\mu) = \max_{q: \mathbb{E}_q[T(x)] = \mu} H(q)$.

As we showed: (1). $A^*(\mu) = \max_{\theta}[\theta^T \mu - A(\theta)]$, $A(\theta) = \max_{\mu}[\theta^T \mu - A^*(\mu)]$

$$\min_{q \in Q} D(q \| p) = A(\theta) - \max_{\mu \in M} [\theta^T \mu + A^*(\mu)]$$

13.6 Connection between Exponential Families and Graphic Models

Pairwise Markov network over $G(V, E)$

$$\begin{aligned} p(\vec{x}) &= \frac{1}{Z} \left(\prod_{i \in V} \psi_i(x_i) \right) \left(\prod_{(i,j) \in E} \psi_{ij}(x_i x_j) \right) \\ &= \frac{1}{Z} e^{\sum_{i \in V} \ln \psi_i(x_i) + \sum_{(i,j) \in E} \ln \psi_{ij}(x_i x_j)} \\ &= \frac{1}{Z} e^{\sum_{i \in V} \sum_{x \in X} \ln \psi_i(x_i) \mathbf{1}_{x=i} + \sum_{(i,j) \in E} \sum_{x, x' \in X} \ln \psi_{ij}(x_i x_j) \mathbf{1}_{x=i, x'=j}} \end{aligned}$$

We can let

$$T_{ix}(x) = \mathbf{1}_{x_i=x}, \quad \forall i \in V, \forall x \in X$$

$$\theta_i(x) = \ln \psi_i(x), \quad \forall i \in V, \forall x \in X$$

$$T_{ijxx'}(x, x') = \mathbf{1}_{x_i=x, x_i=x'}, \quad \forall (i, j) \in E, \forall x, x' \in X$$

$$\theta_{ij}(x) = \ln \psi_{ij}(x, x'), \quad \forall (i, j) \in E, \forall x, x' \in X$$

The probability can be transformed into exponential families.

The dimension of this family is $d = |V||X| + |E||X|^2$.

13.6.1 Marginal polytope

Definition 13.2 (marginal polytope)

The mean parameters associated with the distribution p_θ are the 1-dimensional marginals for the vertices,

$$\mu_i(x) = \mathbb{E}_\theta[T_{ix}(x)] = P_\theta(x_i = x) = \text{marginal distribution of } X_i$$

and the pairwise marginals associated with the edge set E ,

$$\mu_{ij}(x, x') = \mathbb{E}_\theta[T_{ijxx'}(x, x')] = P_\theta(x_i = x, x_j = x') = 2D \text{ marginal distribution of } (X_i, X_j)$$

*The set of realizable mean parameters, \mathcal{M} , is then called the **marginal polytope** and denoted by $\mathcal{M}(G)$.*

13.6.2 Locally Consistent Marginal Distributions

Given a graph $G = (V, E)$, consider the set of marginal distributions τ_i on individual nodes $i \in V$ and pairwise marginals τ_{ij} on edges $(i, j) \in E$ that are locally consistent in the sense

$$\begin{aligned} \sum_x \tau_i(x) &= 1, \quad \tau_i(x) \geq 0, & \forall i \in V, x \in X \\ \sum_{x,x'} \tau_{ij}(x, x') &= 1, \quad \tau_{ij}(x, x') \geq 0, & \forall (i, j) \in E, x, x' \in X \\ \sum_{x_j \in X} \tau_{ij}(x_i, x_j) &= \tau_i(x_i), & \forall (i, j) \in E, x_i, x_j \in X \\ \sum_{x_i \in X} \tau_{ij}(x_i, x_j) &= \tau_j(x_j), & \forall (i, j) \in E, x_i, x_j \in X \end{aligned}$$

Definition 13.3 (local marginal polytope)

The **local marginal polytope** $\mathcal{L}(G)$ is the set of $\tau = (\{\tau_i\}_{i \in V}, \{\tau_{ij}\}_{(i,j) \in E})$ that satisfy the above consistency conditions.



This is a fairly simple polytope defined by $|V| + (2|X| + |X|^2)|E|$ linear constraints. Clearly the marginal polytope $\mathcal{M}(G)$ is a subset of $\mathcal{L}(G)$, but is the converse true?

Proposition 13.1 (For foster $\mathcal{M}(G) = \mathcal{L}(G)$)

If $G = (V, E)$ is a forest then $\mathcal{M}(G) = \mathcal{L}(G)$. Any probability distribution on G can be expressed as follows in terms of its 1-D and pairwise marginals:

$$p(\vec{x}) = \left(\prod_{i \in V} \mu_i(x_i) \right) \left(\prod_{(i,j) \in E} \frac{\mu_{ij}(x_i, x_j)}{\mu_i(x_i) \mu_j(x_j)} \right)$$


Proof 13.1

We first prove the claim for a tree. Any tree can be generated starting from a single node and adding one edge at a time. The claim can be proven by induction. It clearly holds for a graph consisting of two nodes and a single edge (i, j) , since $p(x) = \mu_{ij}(x_i, x_j)$ in this case. If a new node k and a new edge (j, k) are added to an existing tree $(\mathcal{V}', \mathcal{E}')$ where $j \in \mathcal{V}'$, we obtain a new tree $(\mathcal{V}, \mathcal{E})$ with $\mathcal{V} = \mathcal{V}' \cup \{k\}$ and $\mathcal{E} = \mathcal{E}' \cup \{(j, k)\}$. If

$$p(\mathbf{x}_{\mathcal{V}'}) = \left(\prod_{i \in \mathcal{V}'} \mu_i(x_i) \right) \left(\prod_{(i,j) \in \mathcal{E}'} \frac{\mu_{ij}(x_i, x_j)}{\mu_i(x_i) \mu_j(x_j)} \right)$$

then

$$p(\mathbf{x}) = p(\mathbf{x}_{\mathcal{V}'}, x_k) = p(\mathbf{x}_{\mathcal{V}'}) p(x_k | x_j) = p(\mathbf{x}_{\mathcal{V}'}) \frac{\mu_{jk}(x_j, x_k)}{\mu_j(x_j)}$$

satisfies the claim. Next, if the forest contains more than one tree, the distribution $p(\mathbf{x})$ factors over the trees, and the claim still holds.

Finally, to any $\mu \in \mathcal{L}(G)$ we can associate a global distribution p using the claim. The marginals and

pairwise marginals of p are given by μ , hence $\mu \in \mathcal{M}(\mathcal{G})$. This proves the first part of the claim.

Example 13.1 However for a graph that is not a tree, $\mathcal{M}(G)$ is in general a strict subset of $\mathcal{L}(G)$. Consider for instance the 3-cycle with node set $V = \{1, 2, 3\}$ and edge set $E = \{(1; 2); (2; 3); (3; 1)\}$. Let $X = \{0, 1\}$ and consider τ_1, τ_2, τ_3 that are uniform over X , and

$$\tau_{12} = \tau_{23} = \begin{bmatrix} 0.5 - \varepsilon & \varepsilon \\ \varepsilon & 0.5 - \varepsilon \end{bmatrix}, \quad \tau_{31} = \begin{bmatrix} \varepsilon & 0.5 - \varepsilon \\ 0.5 - \varepsilon & \varepsilon \end{bmatrix}$$

for some $\varepsilon \in (0, 0.5)$. By inspection, $\tau \in \mathcal{L}(\mathcal{G})$. However, for ε small enough, the definitions of $\tau_{12}, \tau_{23}, \tau_{31}$ imply respectively that $X_1 = X_2, X_2 = X_3$, and $X_3 \neq X_1$ with high probability. These conditions are incompatible, hence $\tau \notin \mathcal{M}(\mathcal{G})$.

In this example, the edge set is small, and it is relatively easy to determine that $\tau \notin \mathcal{M}(\mathcal{G})$. For a large graph, this would generally not be computationally feasible. Since τ may not be the marginals of any joint distribution on \mathcal{G} , τ are often referred to as **pseudomarginals**.

Definition 13.4 (pseudomarginal)

Marginal distribution τ such that $\tau \in \mathcal{L}(\mathcal{G})$ and $\tau \notin \mathcal{M}(\mathcal{G})$ is called **pseudomarginals**.



13.6.3 Entropy on Tree Graphs

Any distribution p defined on a tree graph is of the form $p(\vec{x}) = (\prod_{i \in V} \mu_i(x_i)) \left(\prod_{(i,j) \in E} \frac{\mu_{ij}(x_i, x_j)}{\mu_i(x_i) \mu_j(x_j)} \right)$. Hence, its entropy is given by

$$\begin{aligned} H(p) &= \mathbb{E}_p[-\ln p(\vec{X})] \\ &= \sum_{i \in V} \mathbb{E}_p[-\ln \mu_i(X_i)] - \sum_{(i,j) \in E} \mathbb{E}_p \left[\ln \frac{\mu_{ij}(X_i, X_j)}{\mu_i(X_i) \mu_j(X_j)} \right] \\ &= \sum_{i \in V} H(\mu_i) - \sum_{(i,j) \in E} I(\mu_{ij}) \end{aligned}$$

where

$$\begin{aligned} I(\mu_{ij}) &\triangleq \mathbb{E}_{\mu_{ij}} \left[\ln \frac{\mu_{ij}(X_i, X_j)}{\mu_i(X_i) \mu_j(X_j)} \right] \\ &= D(\mu_{ij} \| \mu_i \mu_j) \\ &= H(\mu_i) + H(\mu_j) - H(\mu_{ij}) \end{aligned}$$

is the **mutual information** associated with the pairwise marginal μ_{ij} . Since this is a Kullback-Leibler divergence, it is nonnegative. The mutual information is zero if X_i and X_j are independent random variables and is upper-bounded by both $H(\mu_i)$ and $H(\mu_j)$ (value achieved if X_j is a function of X_i , or vice-versa).

The entropy of p is easily computed but is not concave in μ . Equivalently, the set of distributions p on a tree

graph is generally nonconvex. (Consider a length-3 chain for instance.)

13.6.4 Naive Mean-Field Methods In Graph

Approximate complicated p.d.f. $p_\theta(\vec{x}) = e^{\theta^T I(\vec{x}) - A(\theta)}$ in $G = (V, E)$ with tractable $q(\vec{x}) \in G'(V, E')$, where $E' \subset E$. $q(\vec{x}) = \prod_{i \in V} q_i(x_i)$ (naive mean field method assumes fully factorized).

Minimizing divergence:

$$\begin{aligned} D(q \| p_\theta) &= \mathbb{E}_q \left[\ln \frac{q(\vec{x})}{p_\theta(\vec{x})} \right] \\ &= -\theta^T \mathbb{E}_q [T(\vec{x})] + A(\theta) - H(q) \end{aligned}$$

$$Q = \{q : \mathbb{E}_q[T(\vec{x})] = \mu, \mu \in M'\}$$

$$\begin{aligned} \min_{q \in Q} D(q \| p_\theta) &= A(\theta) - \max_{\mu \in M'} [\theta^T \mu - A^*(\mu)] \\ \max_{\{\mu_i\}_{i \in V}} [\theta^T \mu - A^*(\mu)] &= \sum_{i \in V} \sum_{x \in X} \theta_{ix} \mu_i(x) + \sum_{(i,j) \in E} \sum_{x, x'} \theta_{ijxx'} \mu_{ij}(x, x') + \sum_{i \in V} H(\mu_i) \end{aligned}$$

Taking Lagrangian and taking derivative

$$0 = \frac{\partial L(\mu, \lambda)}{\partial \mu_i(x)} \Rightarrow \mu_{i(x)} = \frac{1}{Z} e^{\theta_{ix} + \sum_{i \in N(i)} \sum_{x' \in X} \theta_{ijxx'} \mu_j(x')}$$

13.6.5 Structural Mean Field Optimization

$q(\vec{x}) = q_1(x_1)q_2(x_2|x_1)q_3(x_3|x_2)$ (Markov Chain in a tree $G' = (V, E')$)

$$\mu_{12}(x_1, x_3) = \sum_{x_2} p(x_1, x_2, x_3) = \sum_{x_2} p(x_1, x_2)p(x_3|x_2) = \sum_{x_2} p(x_1, x_2) \frac{p_{23}(x_2, x_3)}{p(x_2)}$$

13.6.6 Bethe Entropy Approximation

When we compute the entropy of a tree graph, the entropy equals to

$$H(p) \triangleq \sum_{i \in V} H(\mu_i) - \sum_{(i,j) \in E} I(\mu_{ij})$$

(only holds in tree graph!).

In a more general situation. For a distribution p that is not defined on a tree graph, $H(p)$ does not admit a simple expression, and cannot be expressed simply in terms of 1-D marginals and pairwise marginals. (Verify on a 3-cycle). However, if these marginals are known, one could use an approximation to $H(p)$.

Definition 13.5 (Bethe approximation)

We use the equation satisfied in tree graph to approximate entropy in general situations. This approxi-

mation is known as the **Bethe approximation**, and the functional

$$H_{\text{Bethe}}(\tau) \triangleq \sum_{i \in V} H(\tau_i) - \sum_{(i,j) \in E} I(\tau_{ij}), \quad \tau \in \mathcal{L}(\mathcal{G})$$

is known as the **Bethe entropy**. This "entropy" is well defined for all pseudomarginals $\tau \in \mathcal{L}(\mathcal{G})$.

The **Bethe variational problem** is defined as

$$A_{\text{Bethe}}(\theta) \triangleq \max_{\tau \in \mathcal{L}(\mathcal{G})} [\theta^\top \tau + H_{\text{Bethe}}(\tau)]$$

and is relatively tractable owing to the simple nature of $\mathcal{L}(\mathcal{G})$ and the availability of a closed-form expression for $H_{\text{Bethe}}(\tau)$.



Compare with the expression

$$A(\theta) = \sup_{\mu \in \mathcal{M}(\mathcal{G})} [\theta^\top \mu + H(p_{\theta(\mu)})]$$

that is unfortunately intractable because of the complex nature of $\mathcal{M}(\mathcal{G})$ and the lack of an explicit form for $H(p_\mu)$. For a general graph, $\mathcal{M}(\mathcal{G}) \subset \mathcal{L}(\mathcal{G})$ and Bethe entropy is an approximation to entropy; $A_{\text{Bethe}}(\theta)$ is not a bound on $A(\theta)$, only an approximation (see example below). For a tree graph however, $\mathcal{M}(\mathcal{G}) = \mathcal{L}(\mathcal{G})$ and $A_{\text{Bethe}}(\theta) = A(\theta)$

Example 13.2 (Inexactness of Bethe approximation) Consider a fully connected graph with four nodes, $V = \{1, 2, 3, 4\}$, uniform 1-D marginals $\mu_i, i \in V$, and pairwise marginals. $\mu_{ij} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}, \forall i, j \in V \Rightarrow X_i = X_j$ w.p.1. $\vec{x} = [0, 0, 0, 0]$ or $[1, 1, 1, 1]$ with probability 0.5 each.

We have $\mu \in \mathcal{M}(\mathcal{G})$; indeed the distribution p that places probability $\frac{1}{2}$ on the sequences $(0, 0, 0, 0)$ and $(1, 1, 1, 1)$ satisfies the marginal constraints above. We have $H(\mu_i) = \ln 2$ for all $i \in \mathcal{V}$ and $I(\mu_{ij}) = \ln 2$ for all $i \neq j \in \mathcal{V}$. Since there are 6 edges, we obtain

$$H_{\text{Bethe}}(\mu) = 4 \ln 2 - 6 \ln 2 = -2 \ln 2 < 0$$

which shows that the Bethe entropy does not satisfy the same properties as an entropy (it can be negative). The actual entropy $H(p) = \ln 2 > 0$.

Chapter 14 ℓ_1 Penalized Least Squares Minimization

We will focus on an ℓ_1 -penalized least-squares problem where the objective function is the sum of a quadratic function representing "fit to the data" and a regularization term which is the ℓ_1 norm of an unknown signal to be recovered. The first term is smooth, the second is not.

14.1 Problem Statement

Given an observation vector $y \in \mathbb{R}^m$, a $m \times n$ matrix A , a constant $\lambda > 0$, find a vector $x \in \mathbb{R}^n$ that achieves the minimum of

$$f(x) \triangleq \frac{1}{2} \|y - Ax\|^2 + \lambda \|x\|_1.$$

The first component is half the squared ℓ_2 norm of $r \triangleq y - Ax$, which can be interpreted as an observation error. The second component, $\|x\|_1 \triangleq \sum_{i=1}^n |x_i|$, is the ℓ_1 norm of x .

The problem admits a Bayesian interpretation, in which the observations y are the sum of Ax and white Gaussian noise with mean zero and variance σ^2 ,

$$y = Ax + z, \text{ where } z \sim \mathcal{N}(0, \sigma^2)$$

and x is a realization of a random vector with iid entries following a double-exponential (Laplace) distribution.

In this case,

$$\ln p(y | x) + \ln p(x) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|y - Ax\|^2 - n \ln 2 - \|x\|_1$$

hence minimizing f is equivalent to MAP estimation for the above Bayesian problem, with $\lambda = 2\sigma^2$.

The structure of A depends on the application. In signal processing and computer vision, A is usually related to a convolution operator, describing for instance motion blur in video. In compressive sensing, the entries of A are typically iid random variables. The algorithms we will focus on do not require A to have any special structure. We begin with two important simple cases (identity and orthonormal A), then move to the general problem.

14.2 Special Cases

14.2.1 Definition: Soft Threshold

Definition 14.1 (Soft Threshold)

$$S_\lambda(y) \triangleq \begin{cases} y - \lambda, & \text{for } y \geq \lambda \\ y + \lambda & \text{for } y \leq -\lambda \\ 0 & \text{for } |y| < \lambda \end{cases}$$



14.2.2 Identity A

Let $m = n$ and A be the identity matrix

$$f(x) \triangleq \sum_{i=1}^n \left(\frac{1}{2}(y_i - x_i)^2 + \lambda|x_i| \right)$$

$$f'(x) = \begin{cases} -y + x + \lambda \text{sign}(x), & x \neq 0 \\ \text{does not exist,} & x = 0 \end{cases}$$

The solution x is obtained by applying a soft threshold to each component y_i of the observations,

$$0 = f'(x) \Rightarrow x = S_\lambda(y) \triangleq \begin{cases} y - \lambda, & \text{for } y \geq \lambda \\ y + \lambda & \text{for } y \leq -\lambda \\ 0 & \text{for } |y| < \lambda \end{cases}$$

14.2.3 Orthonormal A

If $m = n$ and A is orthonormal, then $A^{-1} = A^T$ and

$$\|y - Ax\|^2 = \|A(A^T y - x)\|^2 = \|A^T y - x\|^2$$

$$\Rightarrow x = S_\lambda(A^T y)$$

14.2.4 Quadratic Optimization ($\lambda = 0$)

We now consider general A . It is useful to first study the case $\lambda = 0$, in which case f is quadratic and the optimization problem is smooth. The solution x satisfies the necessary first-order optimality condition

$$0 = \nabla f(x) = -A^\top(y - Ax) \in \mathbb{R}^n.$$

If $\text{rank}(A) \geq n$ (which implies $m \geq n$), the unique solution is $x = (A^\top A)^{-1} A^\top y$.

Otherwise, the solution is nonunique. Any $x = A^+y + z$ where $z \in \text{Null}(A)$ and $A^+ \in \mathbb{R}^{n \times m}$ is the **Moore**

pseudo-inverse of A , is a solution. The minimum-norm solution is $x = A^+y$.

Even though a closed-form solution exists, for large n one would avoid the computationally expensive matrix inverse and use an iterative algorithm such as gradient descent or conjugate gradient to derive the solution. The gradient descent update takes the form

$$x^{k+1} = x^k + \alpha A^\top (y - Ax^k), \quad k = 1, 2, 3, \dots$$

where α is the step size.

14.3 General Solution: Lasso

14.4 General Solution: Iterative Soft Thresholding Algorithm (ISTA)

The idea is to tackle the difficult optimization problem by solving a sequence of simple optimization problems. Often (as is the case here) the simple optimization problems will admit an easily-computable closed-form solution.

14.4.1 Proximal Minimization Algorithm

Definition 14.2 (Proximal Minimization Algorithm)

Consider a convex function F and a $n \times n$ positive definite matrix W . The iterative algorithm with update equation

$$x^{k+1} = \arg \min_{x \in \mathbb{R}^n} \left\{ F(x) + \frac{1}{2} (x - x^k)^T W (x - x^k) \right\}$$

is a **proximal minimization algorithm**.



If W is suitably chosen, the algorithm is a majorization-minimization algorithm.

Lemma 14.1 (Nondecreasing and Convergence)

Such algorithms are of the form

$$x^{k+1} = \arg \min_{x \in \mathbb{R}^n} Q(x, x^k)$$

where the function $Q(\cdot, x')$ is easy to minimize, $Q(x, x) = F(x)$, and $Q(x, x') \geq F(x)$. The sequence $F(x^k)$ is nondecreasing because by application of the above properties and the definition of x^{k+1}

$$F(x^{k+1}) \leq Q(x^{k+1}, x^k) \leq Q(x^k, x^k) = F(x^k).$$

One may also show that the sequence x^k converges to a minimum of F .



14.4.2 Apply to ℓ_1 -penalized least-squares

To apply this strategy to ℓ_1 -penalized least-squares, let $F(x) = \frac{1}{2}\|y - Ax\|^2 + \lambda\|x\|_1$ and

$$Q(x, x') = F(x) + \frac{c}{2}\|x - x'\|^2 - \frac{1}{2}\|\mathbf{A}(x - x')\|^2$$

where c is larger than the maximum eigenvalue of $\mathbf{A}^\top \mathbf{A}$. Then $\mathbf{W} \triangleq c\mathbf{I}_n - \mathbf{A}^\top \mathbf{A}$ is a symmetric positive definite matrix and

$$\begin{aligned} Q(x, x^k) &= F(x) + \frac{c}{2}\|x - x^k\|^2 - \frac{1}{2}\|\mathbf{A}(x - x^k)\|^2 \\ &= F(x) + \frac{1}{2}(x - x^k)^\top \underbrace{(c\mathbf{I}_n - \mathbf{A}^\top \mathbf{A})}_{=\mathbf{W}} (x - x^k) \end{aligned}$$

satisfies the properties of a majorizing function.

We now show that $Q(\cdot, x^k)$ is easy to minimize:

$$\begin{aligned} Q(x, x^k) &= \lambda\|x\|_1 + \frac{1}{2}\|y - Ax\|^2 + \frac{c}{2}\|x - x^k\|^2 - \frac{1}{2}\|\mathbf{A}(x - x^k)\|^2 \\ &= \lambda\|x\|_1 - x^\top \underbrace{\left[\mathbf{A}^\top(y - Ax^k) + cx^k\right]}_{=cu^k} + \frac{c}{2}\|x\|^2 + \text{constant} \\ &= \frac{c}{2}\|x - u^k\|^2 + \lambda\|x\|_1 + \text{constant} \end{aligned}$$

where $u^k \triangleq x^k + \frac{1}{c}\mathbf{A}^\top(y - Ax^k)$. The minimization of $Q(\cdot, x^k)$ takes the same form as *identity A situation*.

Hence, the solution is obtained in closed form using the componentwise soft threshold operator:

$$x^{k+1} = \arg \min_x Q(x, x^k) = S_{\frac{\lambda}{c}}(u^k) = S_{\frac{\lambda}{c}}\left(\frac{1}{c}\mathbf{A}^\top(y - Ax^k) + x^k\right)$$

This is the update equation for the **Iterative Soft Thresholding Algorithm (ISTA)**. Observe that the equation is an extension of the gradient descent update for the purely quadratic problem (with $\lambda = 0$ and step size $\alpha = \frac{1}{c}$).

14.5 Convergence Rate

We say **linear convergence** if

$$\|x^k - x^*\| \leq ab^k$$

where $a > 0$ and $b \in (0, 1)$.

The sufficient condition of **linear convergence** is

$$\lim_{k \rightarrow \infty} \sup \frac{\|x^{k+1} - x^*\|}{\|x^k - x^*\|} \leq \beta$$

for some $\beta \in (0, 1)$.

As an application, we consider $\min_x \{f(x) = x^T Q x\}$ where $Q \succ 0$, achieves x^* at 0.

GD \Rightarrow step size $\alpha < \frac{2}{\lambda_{\max}(Q)}$. The optimal step size $\alpha_{opt} = \frac{2}{\lambda_{\max}(Q) + \lambda_{\min}(Q)}$

Condition number of Q is $\kappa = \frac{\lambda_{\max}(Q)}{\lambda_{\min}(Q)} \geq 1$. With $\alpha = \alpha_{opt}$, $\beta = \frac{\|x^{k+1}\|}{\|x^k\|} = \frac{\kappa-1}{\kappa+1} < 1$

Heavy -Ball Method

$$x^{k+1} = x^k - \alpha \nabla f(x^k) + \beta(x^k - x^{k-1})$$

(with momentum), $\beta > 0$. For minimization of $\min f(x) = x^T Q x$, heavy ball can then be shown to be equivalent to conjugate gradient when α and β are optimized.

14.6 Fast Iterative Soft Thresholding Algorithm (FISTA)

In the ISTA, updates take the form of

$$x^{k+1} = S_{\lambda/c} \left(\frac{1}{c} A^\top (y - Ax^k) + x^k \right)$$

In the FISTA, consider the sequence $t_{k+1} = \frac{1}{2} \left(1 + \sqrt{1 + 4t_k^2} \right)$, initialized with $t_1 = 1$. It can easily be shown that $t_k = \frac{k}{2}[1 + o(1)]$

And the FISTA takes the form:

$$\begin{aligned} x^{k+1} &= S_{\lambda/c} \left(\frac{1}{c} A^\top (y - A\tilde{x}^k) + \tilde{x}^k \right) \\ \tilde{x}^{k+1} &= x^k + \frac{t_k - 1}{t_{k+1}} (x^k - x^{k-1}) \end{aligned}$$

where $\tilde{x}^1 = x^{(0)}$ and the second term in the right side is a "momentum" term, as in the heavy ball algorithm.

The constant c should be larger than the maximum eigenvalue of $A^\top A$.

Analysis of the convergence rate of FISTA-type algorithms is a current research topic.

14.7 Alternating Direction Method of Multipliers (ADMM)

The Alternating Direction Method of Multipliers (ADMM) is an augmented Lagrangian method.

The general idea is to reformulate the minimization problem

$$\min_x \{g(x) + h(x)\}$$

as

$$\min_{x,z} \{g(x) + h(z)\} \quad \text{subj. to} \quad x = z$$

which is solved using an augmented Lagrangian approach.

Identify $g(x) = \frac{1}{2} \|y - Ax\|^2$ and $h(x) = \lambda \|x\|_1$, fix some penalty parameter $\nu > 0$, and write the augmented

Lagrangian as

$$\mathcal{L}(x, z, u) = \frac{1}{2} \|y - Ax\|^2 + \lambda \|z\|_1 + \frac{\nu}{2} \|x - z\|^2 + u^\top (x - z)$$

which is linear in the vector of Lagrange multipliers u , strongly convex in (x, z) , and nonsmooth in z but easy to minimize over z given (x, u) . The update equations are

$$\begin{aligned} x^{k+1} &= \arg \min_x \mathcal{L}(x, z^k, u^k) \\ &= (A^\top A + \nu I)^{-1} (A^\top y + \nu z^k - u^k) \\ z^{k+1} &= \arg \min_z \mathcal{L}(x^{k+1}, z, u^k) \\ &= S_{\lambda/\nu}(x^{k+1} - \nu^{-1} u^k) \\ u^{k+1} &= u^k + \nu (x^{k+1} - z^{k+1}), \quad k = 1, 2, 3, \dots . \end{aligned}$$

Chapter 15 Compressive Sensing

The problem is to recover a sparse signal $x \in \mathbb{R}^n$ by solving $y = Ax$, where $y \in \mathbb{R}^m$, $A \in \mathbb{R}^{m \times n}$, and $x \in \mathbb{R}^n$ is a sparse vector. $m \ll n$. In general this would be a severely underdetermined linear system, but recovery is possible if the signal x is sparse (contains mostly zeroes), or at least approximately sparse.

15.1 Definitions related to Sparsity

Definitions:

1.

Definition 15.1 (ℓ_p norm)

The ℓ_p norm of a vector $x \in \mathbb{R}^n$ is

$$\|x\|_p \triangleq \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} \text{ for all } p \geq 1$$



2.

Definition 15.2 (ℓ_0 pseudonorm)

The ℓ_0 pseudonorm of a vector $x \in \mathbb{R}^n$ is

$$\|x\|_0 \triangleq \sum_{i=1}^n \mathbf{1}_{\{x_i \neq 0\}}$$

i.e., the number of nonzero components of x .



3.

Definition 15.3 (k -sparse)

A signal $x \in \mathbb{R}^n$ is k -sparse if

$$\|x\|_0 \leq k$$

i.e., the number of nonzero components of x is smaller than k .



4.

Definition 15.4 (set of k -sparse signals)

The set of k -sparse signals is

$$\Sigma_k \triangleq \{x \in \mathbb{R}^n : \|x\|_0 \leq k\}$$



This is a union of $\binom{n}{k}$ -dimensional subspaces of \mathbb{R}^n , which is not a linear space.

For instance let $n = 2$ and $k = 1$, then Σ_1 is the union of the horizontal and vertical axes in the 2-D plane.

Many signals are sparse in a transform domain (for instance, Fourier-sparse) or approximately sparse.

For instance, the wavelet coefficients of an image or a speech signal are approximately sparse, and one can typically construct good approximations of such signals by using only the largest 5% of their components.

Claim 15.1 (Producing k -sparse signal)

*One can go from an approximately sparse to an (exactly) k -sparse signal by applying a **hard threshold operator**, $\hat{x} = H_k(x)$, producing a vector in which all but the k largest (in magnitude) coefficients of x are set to zero.*



5.

Definition 15.5 (ℓ_p approximation error)

The ℓ_p approximation error of $x \in \mathbb{R}^n$ in Σ_k is

$$e_{k,p}(x) \triangleq \min_{\hat{x} \in \Sigma_k} \|x - \hat{x}\|_p, \quad p \geq 1$$



Lemma 15.1 (Properties)

(a). If x is k -sparse, i.e., $x \in \Sigma_k$, then $e_{k,p}(x) = 0$ for all p . (Optimal \hat{x} is exactly the x)

(b). Otherwise, the approximation error typically vanishes geometrically with k , i.e.,

$$e_{k,p}(x) \leq ck^{-r}$$

for some $c, r > 0$, generally dependent on x .

(c). It is easily shown that the minimum of error is achieved by $\hat{x} = H_k(x)$.



Other forms of sparsity are frequently encountered:

- **Low-dimensional manifolds:** let Θ be a compact subset of \mathbb{R}^k and $f : \Theta \rightarrow \mathbb{R}^n$ a continuously differentiable mapping, then $x = f(\theta), \theta \in \Theta$ belongs to a k -dimensional manifold embedded in \mathbb{R}^n . This model applies to face images under varying illumination.
- **Low-rank matrices:** let x be a $n_1 \times n_2$ matrix of rank $r < \min(n_1, n_2)$. Then x can be represented using the singular value decomposition $x = \sum_{j=1}^r \sigma_j u_j v_j^\top$ where σ_j are the (positive) singular values, and $u_j \in \mathbb{R}^{n_1}$ and $v_j \in \mathbb{R}^{n_2}$ are the singular vectors for $1 \leq j \leq r$. This model finds applications to computer vision, geolocalization, and collaborative filtering (cf the Netflix recommender system).

15.2 Measurement Matrix

Definition 15.6 (Measurement Matrix)

The measurement matrix A is a $m \times n$ matrix where $m \ll n$ (fat matrix). The observations are given by

$$\vec{y} = A\vec{x} \in \mathbb{R}^m \text{ in case of noise-free measurements.}$$



In this section we consider two basic questions are:

1. What properties should $A \in \mathbb{R}^{m \times n}$ have so that $\vec{x} \in \mathbb{R}^n$ can be recovered from $\vec{y} \in \mathbb{R}^m$?
2. If A satisfies those properties, what recovery algorithm can be used?

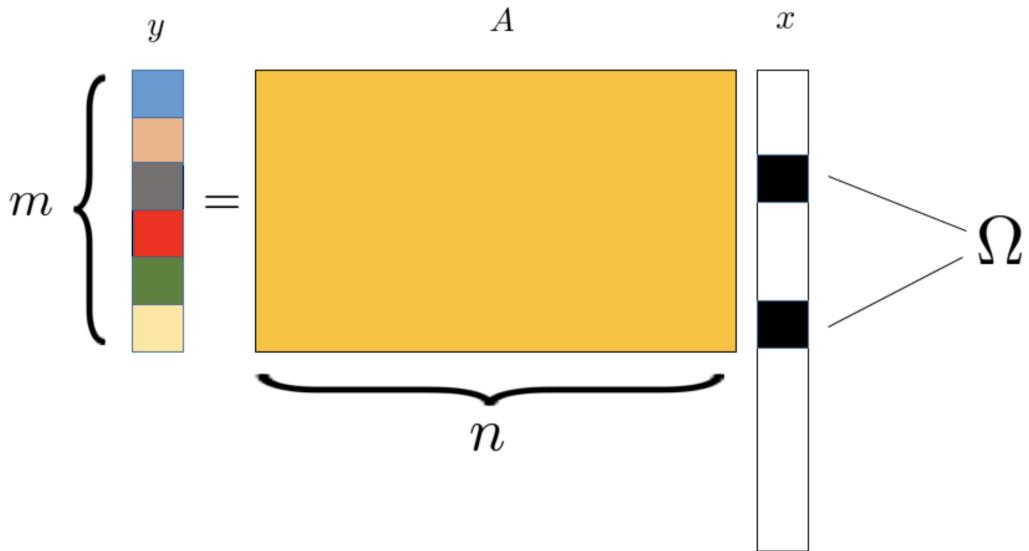


Figure 15.1: Measurement of sparse signal \vec{x} with support set Ω of size k .

Denote by A_i , $1 \leq i \leq n$ the column vectors of A and let $\Omega \subset \{1, 2, \dots, n\}$ be the support set of the vector \vec{x} , i.e., $\Omega = \{i : x_i \neq 0\}$. Define the reduced matrix $A_\Omega = \{A_i, i \in \Omega\}$ and the reduced vector $\vec{x}_\Omega = \{x_i, i \in \Omega\}$. Then

$$\vec{y} = A\vec{x} = A_\Omega\vec{x}_\Omega$$

15.2.1 Matrix Preliminaries

The SVD decomposition of A takes the form $A = U[D \mid 0]V^T$ where U and V are $m \times m$ and $n \times n$ unitary matrices (whose conjugate transpose equals to inverse), respectively, $D = \text{diag}(\sigma_k)_{k=1}^m$, and 0 is the $m \times (n - m)$ all-zero matrix. The m nonnegative entries σ_k , $1 \leq k \leq m$, are the singular values of A . The

matrix A can be expanded as a sum of m rank-one matrices, $A = \sum_{k=1}^m \sigma_k u_k v_k^\top$ where $\{u_k\}$ and $\{v_k\}$ are the columns of U and V , respectively.

The spectral norm of A is

$$\|A\| \triangleq \sup_{\vec{x}: \|\vec{x}\|_2=1} \|A\vec{x}\|_2$$

and is equal to the largest singular value of A .

The $m \times m$ Gram matrix $G \triangleq AA^\top$ is symmetric and nonnegative definite and can be expressed as $G = UD^2U^\top$ hence its eigenvalues are given by $\lambda_k(G) = \sigma_k^2(A)$, $1 \leq k \leq m$.

Definition 15.7 (Null Space)

The null space of A is

$$N(A) \triangleq \{\vec{x} \in \mathbb{R}^n : A\vec{x} = 0\}$$

This is a subspace of \mathbb{R}^n , whose dimension is at most $n - m$.



Definition 15.8 (Spark)

The spark of A is the smallest number of columns of A that are linearly dependent.

If any set of q columns of the matrix are linearly independent, $\text{spark}(A) = q + 1 \in [2, m + 1]$.



The minimum value of the spark of any matrix is 2, and is achieved by any matrix that has two identical columns. Computing the spark of a matrix is an NP-hard problem. In contrast, the rank of a matrix, which is the largest number of columns that are linearly independent, is easy to compute.

15.2.2 Recovery of k -Sparse Signals

Lemma 15.2

Unique recovery of $x \in \Sigma_k$ given measurement matrix A is possible $\Rightarrow \Sigma_{2k} \cap N(A) = \emptyset$



Proof 15.1

To recover any k -sparse signal, we need that $A\vec{x} \neq A\vec{x}'$ for any distinct $\vec{x}, \vec{x}' \in \Sigma_k$. Hence, $\vec{x} - \vec{x}' \notin N(A)$. Since $\vec{x} - \vec{x}'$ is a $2k$ -sparse signal and in fact $\Sigma_{2k} = \vec{x} - \vec{x}' : \vec{x}, \vec{x}' \in \Sigma_k$, we need $\Sigma_{2k} \cap N(A) = \emptyset$.

Theorem 15.1

Unique recovery of $\vec{x} \in \Sigma_k$ given measurement matrix A is possible $\Leftrightarrow \text{spark}(A) > 2k$. ($q \geq 2k$)



Proof 15.2

Eldar, Y. C., & Kutyniok, G. (Eds.). (2012). Compressed sensing: theory and applications, Chapter 1.

15.2.3 Restricted Isometry Property

Definition 15.9 (Restricted Isometry Property (RIP))

The matrix A satisfies the RIP property of order k if there exists $\delta_k \in (0, 1)$ such that

$$(1 - \delta_k) \|\vec{x}\|_2^2 \leq \|A\vec{x}\|_2^2 \leq (1 + \delta_k) \|\vec{x}\|_2^2, \quad \forall \vec{x} \in \Sigma_k$$



If A satisfies the RIP property of order k , then A approximately preserves the ℓ_2 distance between any pair $\vec{x}, \vec{x}' \in \Sigma_k$. This provides a *stable embedding* of k -sparse signals in \mathbb{R}^m . This property will be key to derive a recovery algorithm that is robust to noise.

Claim 15.2 (Equivalent Formulation of RIP)

Since $\vec{y} = A\vec{x} = A_\Omega \vec{x}_\Omega$, an **equivalent formulation** of the RIP is

$$(1 - \delta_k) \|\vec{x}_\Omega\|_2^2 \leq \|A_\Omega \vec{x}_\Omega\|_2^2 \leq (1 + \delta_k) \|\vec{x}_\Omega\|_2^2, \quad \forall \Omega : |\Omega| = k$$

for all Ω of size k and for all $\vec{x}_\Omega \in \mathbb{R}^k$.

$$\begin{aligned} \frac{\|A_\Omega \vec{x}_\Omega\|_2^2}{\|\vec{x}_\Omega\|_2^2} - 1 &= \frac{\vec{x}_\Omega^T (A_\Omega^T A_\Omega - I_k) \vec{x}_\Omega}{\vec{x}_\Omega^T \vec{x}_\Omega} \in [-\delta_k, \delta_k], \quad \forall \Omega : |\Omega| = k \\ \Rightarrow \delta_k &= \max_{\Omega : |\Omega|=k} \|A_\Omega^T A_\Omega - I_k\| \end{aligned}$$



Theorem 15.2 (Measurement Bound)

Assume $A \in \mathbb{R}^{m \times n}$ satisfies the RIP of order $2k$ with RIP constant $\delta_{2k} \in (0, \frac{1}{2}]$. Then

$$m \geq ck \ln \frac{n}{k}$$

where the constant $c = \frac{1}{2} \ln(1 + \sqrt{24}) \approx 0.28$.



15.3 Robust Signal Recovery from Noiseless Observations

We now consider the so-called robust CS problem: the signal \vec{x} is approximately sparse, and the observations are noiseless. A reasonable attempt to recover a sparse signal would be the so-called ℓ_0 recovery problem

$$\min_{\vec{x} \in \mathbb{R}^n} \|\vec{x}\|_0 \quad \text{subj. to} \quad A\vec{x} = \vec{y}$$

Unfortunately this problem is highly nonconvex. Solving it essentially requires evaluating all possible support sets of \vec{x} , which is a combinatorial problem.

A reasonable substitute is the so-called ℓ_1 recovery problem

$$\min_{\vec{x} \in \mathbb{R}^n} \|\vec{x}\|_1 \quad \text{subj. to} \quad A\vec{x} = \vec{y}$$

This procedure tends to produce sparse solutions, as illustrated in the Figure: the line $y = Ax$ is tangent to the

ℓ_1 ball $\|\mathbf{x}\|_1 = \text{cst}$ at the solution \mathbf{x} , producing a 1-sparse solution.

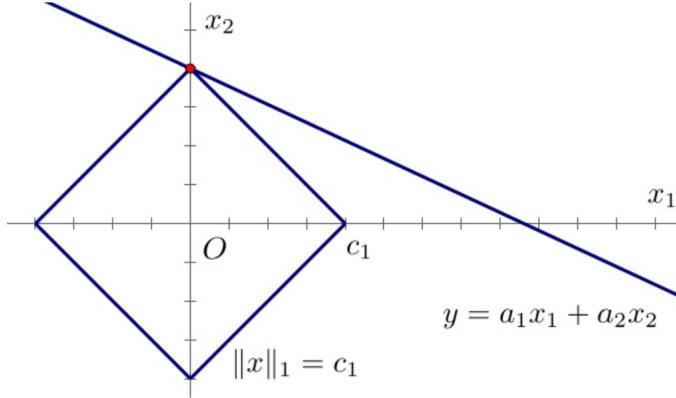


Figure 15.2: ℓ_1 recovery for $n = 2$ and $m = 1$

The following fundamental theorem shows that the ℓ_1 recovery procedure is remarkably good.

Theorem 15.3 (ℓ_1 recovery procedure is good)

Assume A satisfies the RIP of order $2k$ with constant $\delta_{2k} < \sqrt{2} - 1$. $\hat{x} = \operatorname{argmin}_{\vec{x} \in \mathbb{R}^n : A\vec{x} = \vec{y}} \|\vec{x}\|_1$. Then

$$\|\hat{x} - x\|_2 \leq \frac{c}{\sqrt{k}} e_{k,1}(x)$$

and

$$\|\hat{x} - x\|_1 \leq c e_{k,1}(x)$$

with $c = 2 \frac{1-(1-\sqrt{2})\delta_{2k}}{1-(1+\sqrt{2})\delta_{2k}}$. $e_{k,1}(x) \triangleq \min_{\hat{x} \in \Sigma_k} \|x - \hat{x}\|_1$ is the ℓ_1 approximation error of x in Σ_k .



Proof 15.3

Based on the triangle inequality and the inequality $\frac{\|u\|_1}{\sqrt{k}} \leq \|u\|_2 \leq \sqrt{k}\|u\|_\infty$ for all $u \in \Sigma_k$.

Corollary 15.1

If $x \in \Sigma_k$ then $\hat{x} = x$ (exact recovery).



If $x \notin \Sigma_k$ then the quality of the reconstruction is nearly as good as if an *oracle* gave us the location of the k largest absolute components and we measured those directly. (The oracle produces $\hat{x} = H_k(x)$, achieving $e_{k,p}(x)$ for all $p \geq 1$.)

Since $\delta_{2k} < \sqrt{2} - 1 < \frac{1}{2}$, the measurement bound $m \geq ck \ln \frac{n}{k}$ applies, and we need as few as $m = O(k \ln \frac{n}{k})$ measurements to satisfy the conditions of the theorem.

15.4 Robust Signal Recovery from Noisy Observations

15.4.1 Bounded Noise

We consider observations corrupted by bounded noise: $y = Ax + z$ where $\|z\|_2 \leq \epsilon$. We study the ℓ_1 recovery problem

$$\min_{x \in \mathbb{R}^n} \|x\|_1 \quad \text{subj. to} \quad \|Ax - y\|_2 \leq \epsilon$$

which is closely related to the Lasso problem and can be solved using the algorithms introduced in the previous chapter.

Theorem 15.4

Assume A satisfies the RIP of order $2k$ with constant $\delta_{2k} < \sqrt{2} - 1$. Then

$$\|\hat{x} - x\|_2 \leq \frac{c_0}{\sqrt{k}} e_{1,k}(x) + c_1 \epsilon$$

with constants

$$c_0 = 2 \frac{1 - (1 - \sqrt{2})\delta_{2k}}{1 - (1 + \sqrt{2})\delta_{2k}}, \quad \text{and} \quad c_1 = 4 \frac{\sqrt{1 + \delta_{2k}}}{1 - (1 + \sqrt{2})\delta_{2k}}.$$



For $\delta_{2k} = \frac{1}{4}$ the theorem holds with $c_0 \leq 5.5$ and $c_1 \leq 6$. For $\epsilon = 0$, the result coincides with that given in the previous section (noise-free case).