



# Statistical Inference

**Author:** Wenxiao Yang

**Institute:** Haas School of Business, University of California Berkeley

**Date:** 2023

*All models are wrong, but some are useful.*

# Contents

<b>Chapter 1 Statistics Basics</b>	<b>1</b>
1.1 Random Sampling . . . . .	1
1.1.1 Sample Mean and Sample Variance . . . . .	1
1.1.2 Distributional Properties . . . . .	1
1.1.3 Order Statistics . . . . .	2
1.2 Basic Statistics . . . . .	2
1.3 Point Estimation . . . . .	3
1.3.1 Method of Moments (MM) . . . . .	3
1.3.2 Maximum Likelihood (ML) . . . . .	5
1.3.3 Comparing Estimators: Mean Squared Error . . . . .	6
1.3.4 Sufficiency . . . . .	7
<b>Chapter 2 Information-Theoretic Functional</b>	<b>8</b>
2.1 Definitions . . . . .	8
2.1.1 Entropy . . . . .	8
2.1.2 Kullback-Leibler Divergence . . . . .	8
2.1.3 Cross-Entropy . . . . .	9
2.1.4 Mutual Information . . . . .	9
<b>Chapter 3 Statistical Inference</b>	<b>10</b>
3.1 Basics . . . . .	10
3.2 Decision Rule Examples . . . . .	10
3.3 Maximum-Likelihood Principle (state is norandom) . . . . .	11
3.4 Bayesian Decision Rule (state is random) . . . . .	12
3.4.1 Rules . . . . .	12
3.4.2 Maximum A Posteriori (MAP) Decision Rule (Binary example) . . . . .	13
3.4.3 Minimum Mean Squared Error (MMSE) Rule ( $\mathbb{R}^n$ example) . . . . .	13
3.5 Comparison . . . . .	14
<b>Chapter 4 Machine Learning in Inference</b>	<b>15</b>

4.1	Empirical Risk Minimization (ERM) . . . . .	15
4.1.1	Example: Linear MMSE (LMMSE) estimator . . . . .	15
4.1.2	Penalized ERM . . . . .	16
4.2	Stochastic Approximation . . . . .	17
4.3	Stochastic Gradient Descent (SGD) . . . . .	19
4.4	SGD Application to Empirical Risk Minimization (ERM) . . . . .	20
4.4.1	Different Gradient Descent for ERM . . . . .	21
4.4.2	Constraints on Learning Problem . . . . .	21
<b>Chapter 5</b>	<b>Stochastic Integration Methods</b>	<b>23</b>
5.1	Deterministic Methods (Better in Low Dimension) . . . . .	23
5.1.1	Riemann Integration . . . . .	23
5.1.2	Trapezoidal Rule . . . . .	23
5.1.3	Multidimensional Integration . . . . .	24
5.2	Stochastic Methods (Better in High Dimension) . . . . .	24
5.2.1	Classical Monte Carlo Integration . . . . .	24
5.2.2	Importance Sampling . . . . .	25
<b>Chapter 6</b>	<b>Bootstrap (not enough data)</b>	<b>27</b>
6.1	Residual Bootstrap . . . . .	27
<b>Chapter 7</b>	<b>Particle Filtering</b>	<b>29</b>
7.1	Kalman Filtering (Linear Dynamic System) . . . . .	29
7.2	Particle Filtering (Nonlinear Dynamic System) . . . . .	29
7.2.1	Bayesian Recursive Filtering . . . . .	30
7.2.2	Particle Filter (bootstrap filter) . . . . .	30
<b>Chapter 8</b>	<b>EM Algorithm</b>	<b>32</b>
8.1	General Structure of the EM Algorithm . . . . .	32
8.2	Example 1: Variance Estimation . . . . .	34
8.2.1	Maximum-Likelihood (ML) Estimation . . . . .	34
8.2.2	EM Algorithm . . . . .	34
8.3	Example 2: Estimation of Gaussian Mixtures . . . . .	35
8.3.1	Unknown Means: ML estimation is hard . . . . .	35

---

8.3.2 Unknown Means: EM Algorithm . . . . .	36
8.3.3 Unknown Mixture Probabilities, Means and Variances . . . . .	37
8.4 Convergence of EM Algorithm . . . . .	37
8.5 EM As an Alternating Maximization Algorithm . . . . .	38
<b>Chapter 9 Hidden Markov model (HMM)</b>	<b>40</b>
9.1 Viterbi Algorithm: (MAP) estimate $X_{1:t}$ given $Y_{1:t}$ . . . . .	40
9.1.1 MAP estimation problem . . . . .	40
9.1.2 Viterbi Algorithm . . . . .	41
9.2 Bayesian Estimation of a Sequence: Need (MMSE) estimate $X_{1:t}$ given $Y_{1:t}$ . . . . .	42
9.3 Forward-Backward Algorithm: (MMSE) estimate $X_{1:t+1}$ given $Y_{1:t}$ . . . . .	42
9.3.1 $\gamma_t(x) \triangleq P\left\{X_t = x \mid \vec{Y} = \vec{y}\right\}$ . . . . .	42
9.3.2 $\xi_t(x, x') \triangleq P\left\{X_t = x, X_{t+1} = x' \mid \vec{Y} = \vec{y}\right\}$ . . . . .	44
9.3.3 Scaling Factors . . . . .	44
<b>Chapter 10 Graphic Models</b>	<b>45</b>
10.1 Graph Theory . . . . .	45
10.2 Bayesian Networks . . . . .	46
10.3 Markov Networks . . . . .	46
10.3.1 General Form . . . . .	46
10.3.2 Hammersley-Clifford theorem . . . . .	47
10.3.3 Form of Gibbs distribution (Boltzmann distribution) . . . . .	47
10.4 Conversion of directed graph to undirected graph . . . . .	48
10.5 Inference and Learning . . . . .	48
10.5.1 Inference on Trees . . . . .	48
<b>Chapter 11 Variational Inference, Mean-Field Techniques</b>	<b>51</b>
11.1 Naive Mean-Field Methods . . . . .	51
11.1.1 Graphical Models . . . . .	52
11.1.2 Ising Model . . . . .	52
11.2 Exponential Families of Probability Distributions . . . . .	54
11.3 ML Estimation . . . . .	56
11.4 Maximum Entropy . . . . .	56
11.5 . . . . .	57

---

11.6 Connection between Exponential Families and Graphic Models . . . . .	58
11.6.1 Marginal polytope . . . . .	58
11.6.2 Locally Consistent Marginal Distributions . . . . .	58
11.6.3 Entropy on Tree Graphs . . . . .	60
11.6.4 Naive Mean-Field Methods In Graph . . . . .	61
11.6.5 Structural Mean Field Optimization . . . . .	61
11.6.6 Bethe Entropy Approximation . . . . .	61
<b>Chapter 12 <math>\ell_1</math> Penalized Least Squares Minimization</b>	<b>63</b>
12.1 Problem Statement . . . . .	63
12.2 Special Cases . . . . .	64
12.2.1 Definition: Soft Threshold . . . . .	64
12.2.2 Identity $A$ . . . . .	64
12.2.3 Orthonormal $A$ . . . . .	64
12.2.4 Quadratic Optimization ( $\lambda = 0$ ) . . . . .	64
12.3 General Solution: Lasso . . . . .	65
12.4 General Solution: Iterative Soft Thresholding Algorithm (ISTA) . . . . .	65
12.4.1 Proximal Minimization Algorithm . . . . .	65
12.4.2 Apply to $\ell_1$ -penalized least-squares . . . . .	66
12.5 Convergence Rate . . . . .	66
12.6 Fast Iterative Soft Thresholding Algorithm (FISTA) . . . . .	67
12.7 Alternating Direction Method of Multipliers (ADMM) . . . . .	67
<b>Chapter 13 Compressive Sensing</b>	<b>69</b>
13.1 Definitions related to Sparsity . . . . .	69
13.2 Measurement Matrix . . . . .	71
13.2.1 Matrix Preliminaries . . . . .	71
13.2.2 Recovery of k-Sparse Signals . . . . .	72
13.2.3 Restricted Isometry Property . . . . .	73
13.3 Robust Signal Recovery from Noiseless Observations . . . . .	73
13.4 Robust Signal Recovery from Noisy Observations . . . . .	75
13.4.1 Bounded Noise . . . . .	75

# Chapter 1 Statistics Basics

## 1.1 Random Sampling

### Definition 1.1 (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 1.2 (Statistic)

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**.



### 1.1.1 Sample Mean and Sample Variance

#### Definition 1.3 (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  $\mu$  and variance  $\sigma^2$ .

#### Theorem 1.1 ( $\mathbb{E}(\bar{X}), \text{Var}(\bar{X}), \mathbb{E}(S^2)$ )

Suppose  $X_1, \dots, X_n$  is a random sample from a distribution with mean  $\mu$  and variance  $\sigma^2$  (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).  $\mathbb{E}(S^2) = \sigma^2$ .



### 1.1.2 Distributional Properties

#### Theorem 1.2

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

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



**Theorem 1.3 ("Asymptotics")**

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

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



### 1.1.3 Order Statistics

**Definition 1.4 (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 1.1 (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

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



## 1.2 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  $\theta$ ?";
- Hypothesis testing  $\Rightarrow$  "Is  $\theta = \theta_0$ ?";
- Interval estimation  $\Rightarrow$  "Which values of  $\theta$  are 'plausible'?".

**Example 1.1 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.

## 1.3 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 1.5 (Point Estimator)

A **point estimator** (of  $\theta$ ) 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.

### 1.3.1 Method of Moments (MM)

#### Definition 1.6 (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 1.2

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  $\theta$ . Concern  $X_i > \hat{\theta}_{MM}$  could happen. So,  $\max\{\hat{\theta}_{MM}, X_{(n)}\}$  can be better.

### Definition 1.7 (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 1.3

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  $\mu$ ; An alternative better estimator of  $\sigma^2$  is  $\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ .

### 1.3.2 Maximum Likelihood (ML)

#### Definition 1.8 (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

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



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

#### Definition 1.9 (Log-Likelihood)

The **log-likelihood** function is

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



#### Example 1.4

- 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

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

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

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

$$\Rightarrow \hat{\theta}_{ML} = X_{(n)}$$

**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,

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

### 1.3.3 Comparing Estimators: Mean Squared Error

#### General Approach

- o Statistical Decision Theory

Leading Special Case: Mean Squared Error.

#### Definition 1.10 (Mean Squared Error)

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

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

#### Definition 1.11 (Bias)

The **bias** of  $\hat{\theta}$  is (the function of  $\theta$ ) 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 1.12 (Optimality)**

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



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

**1.3.4 Sufficiency****Definition 1.13 (Sufficient Statistic)**

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

**Theorem 1.4 (Rao-Blackwell Theorem)**

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

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

**Proof 1.1**

- (a). Estimator:  $\hat{\theta} = \mathbb{E}[\tilde{\theta} | T]$  doesn't depend on  $\theta$  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  $\text{Var}(\hat{\theta}|T) = 0$  (which also makes  $\hat{\theta} = \tilde{\theta}$ ).

# Chapter 2 Information-Theoretic Functional

## 2.1 Definitions

### 2.1.1 Entropy

#### Definition 2.1 (Entropy)

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

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

(concave in  $p$ )



### 2.1.2 Kullback-Leibler Divergence

#### Definition 2.2 (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 2.1 (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 2.1 (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 2.2 (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  $\frac{q(x)}{p(x)} = 1$ .

### 2.1.3 Cross-Entropy

#### Definition 2.3 (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  $q = p$ .



### 2.1.4 Mutual Information

#### Definition 2.4 (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 3 Statistical Inference

### 3.1 Basics

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

**Example:**

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

### 3.2 Decision Rule Examples

#### Binary HT Example

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

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

Type-I error probability: probability that  $\theta$  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  $\theta$  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  $\tau$ .  $Q(t) = \int_t^\infty \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}} dx$

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

#### Multiple HT Example

Consider three state  $S = \{1, 2, 3\}$ . We can find a  $\tau$  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  $\theta$  is  $i$  but receive  $\delta(x) = j$  (6 types in this example)

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

## Estimation Example

Ex:  $P_\theta \sim N(\theta, \sigma^2)$ . Perform  $\delta(x) = \hat{\theta}$  by using mean-squared error (MSE):

$$MSE = \mathbb{E}_\theta [(\delta(x) - \theta)^2], \theta \in \mathbb{R}$$

## 3.3 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:  $P_0 \sim N(0, \sigma^2), P_1 \sim N(\mu, \sigma^2)$ .

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

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  $H_1$ ,  $\bar{x} \sim N(0, \frac{\sigma^2}{n})$ .

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

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

## 3.4 Bayesian Decision Rule (state is random)

### 3.4.1 Rules

Prior probability distribution  $\pi(\theta)$ ,

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  $l(a, \theta) = |a - \theta|$ .

**Risk of decision rule  $\delta$ :**

$$R(\delta) = \mathbb{E}(l(\delta(X), \theta))$$

where  $(X, \theta)$  are random with prob  $\pi(\theta), P_\theta$

**Note:** to help be consistent with machine learning notations, we use  $y$  to substitute  $\theta$ .

The joint probability  $P(x, y) = \pi(y)P_y(x) = P(x)\pi(y|x)$ .

**Example 3.1** The risk of decision  $\delta(x)$  in Hamming/zero-one loss  $l(a = \hat{y}, y) = \mathbf{1}_{a \neq y}$

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

**Bayes rule**

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

Derive Bayes rule

$$\begin{aligned} R(\delta) &= \int_x \int_y P(x, y) l(\delta(x), y) dy dx \\ &= \int_x P(x) \int_y \pi(y|x) l(\delta(x), y) dy dx \end{aligned}$$

Solve optimization problem:

$$\min_{\delta} \int_x P(x) \int_y \pi(y|x) l(\delta(x), y) dy dx$$

this problem can be transformed into optimization problems for each  $x \in S$

$$\min_{\delta(x)} \int_y \pi(y|x) l(\delta(x), y) dy$$

The problem becomes to compute  $\pi(y|x)$ . From  $P(x, y) = \pi(y)P_y(x) = P(x)\pi(y|x)$ , we know

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

### 3.4.2 Maximum A Posteriori (MAP) Decision Rule (Binary example)

**Example 3.2** 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

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

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

**Example 3.3** (estimation) Squared error loss  $l(a, y) = (a - y)^2$ .

### Minimum Mean Squared Error (MMSE) Rule:

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

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

$$\Rightarrow \delta_B(x) = \mathbb{E}[Y|X = x]$$

which is called **conditional mean estimation**.

In this rule the whole optimization problem also goes to

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

**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.

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

and the posterior risk is independent of  $x$ :

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

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

## 3.5 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 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.

### 4.1 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  $n \rightarrow \infty$ .

#### 4.1.1 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 4.1

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



#### Definition 4.2

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



#### Definition 4.3

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



#### Definition 4.4

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



#### Definition 4.5

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  $a \rightarrow \infty$ )



### 4.1.2 Penalized ERM

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

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

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

Approach 1: Fix  $d << n$ , use ERM.

Approach 2: **(Penalized ERM)**

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

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

## 4.2 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 4.1**  $Y = X + Z$  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 4.1

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

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

- 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 4.2** 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  $\beta$  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 4.3** 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  $a_n = n^{-\alpha}, \alpha \in (\frac{1}{2}, 1] \Rightarrow \beta = \frac{1}{6}, \alpha = \frac{2}{3} \Rightarrow e_n \sim O(n^{-\frac{1}{6}})$

### 4.3 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 4.4**  $f(x) = \frac{1}{2}x^2, x \in \mathbb{R}$ . 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.)

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

#### 4.4.1 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  $\rho$  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  $O(\frac{d}{\varepsilon})$ .

#### 4.4.2 Constraints on Learning Problem

Why achieving a low value of  $\rho$  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  $\delta_B$ ) 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 5 Stochastic Integration Methods

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

## 5.1 Deterministic Methods (Better in Low Dimension)

### 5.1.1 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  $n \sim O(\varepsilon^{-1})$

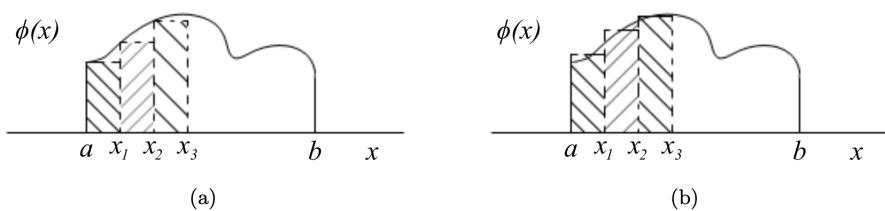
### 5.1.2 Trapezoidal Rule

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



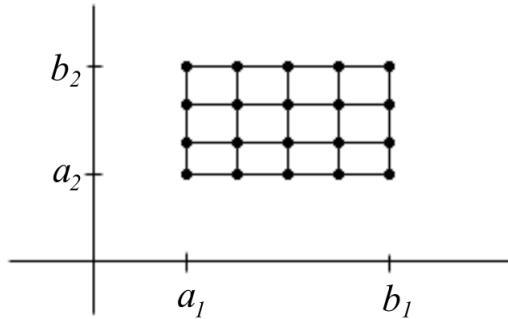
**Figure 5.1:** (a) Riemann approximation; (b) Trapezoidal approximation.

### 5.1.3 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  $\varepsilon$ .

This phenomenon is known as the curse of dimensionality.



**Figure 5.2:** Two-dimensional integration using regular grid.

## 5.2 Stochastic Methods (Better in High Dimension)

### 5.2.1 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  $\xi$  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  $d < 4$ .

### 5.2.2 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  $\xi$

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

$\hat{\xi}_n \xrightarrow{a.s.} \xi$  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  $\delta$  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*( $\xi$ ). 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  $\delta$  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  $\frac{p(x)}{q(x)} h(x) = \alpha, \forall x \in X$ ,  $\alpha$  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 6 Bootstrap (not enough data)

Problem: analyze the performance of an estimator  $\hat{\theta}_n(\vec{Y})$ ,  $\vec{Y} = (Y_1, Y_2, \dots, Y_n)$  taken i.i.d. from distribution  $P$ .

e.g.  $P_\theta = N(0, 1)$ ,  $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n Y_i$

Assume  $\theta$  is a scalar parameter. Performance: (1) Bias  $\mathbb{E}_\theta[\hat{\theta}_n(\vec{Y})] - \theta$ ; (2) Variance  $\mathbb{E}_\theta[\hat{\theta}_n^2(\vec{Y})] - \mathbb{E}_\theta^2[\hat{\theta}_n(\vec{Y})]$ ;

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

### Approach #1 Monte-Carlo Simulations

Generate  $k$  vectors  $\vec{Y}^{(i)}$ ,  $i = 1, 2, \dots, k$  (total  $kn$  random variables) (1) Bias  $\frac{1}{k} \sum_{j=1}^k \hat{\theta}_n(\vec{Y}^{(j)}) - \theta$ ; (2) Variance  $\frac{1}{k} \sum_{j=1}^k \hat{\theta}_n^2(\vec{Y}^{(j)}) - \left(\frac{1}{k} \sum_{j=1}^k \hat{\theta}_n(\vec{Y}^{(j)})\right)^2$ ; (3) CDF  $\hat{G}_n(t) = \frac{1}{k} \sum_{j=1}^k \mathbf{1}_{\hat{\theta}_n(\vec{Y}^{(j)}) < t}$ ,  $\forall t$

### Approach #2 Bootstrap

(When data is not enough) Suppose we only have one data  $\vec{Y} = (Y_1, \dots, Y_n)$

Reuse  $Y_1, \dots, Y_n$  to obtain resamples  $\vec{Y}^* = (Y_1^*, \dots, Y_n^*)$ . Do this  $k$  times  $\Rightarrow k$  resamples  $\vec{Y}^{*(1)}, \dots, \vec{Y}^{*(k)}$

(1) Bias  $\frac{1}{k} \sum_{j=1}^k \hat{\theta}_n(\vec{Y}^{*(j)}) - \theta$ ; (2) Variance  $\frac{1}{k} \sum_{j=1}^k \hat{\theta}_n^2(\vec{Y}^{*(j)}) - \left(\frac{1}{k} \sum_{j=1}^k \hat{\theta}_n(\vec{Y}^{*(j)})\right)^2$ ; (3) CDF  $\hat{G}_n(t) = \frac{1}{k} \sum_{j=1}^k \mathbf{1}_{\hat{\theta}_n(\vec{Y}^{*(j)}) < t}$ ,  $\forall t$

**Example:**  $\theta = \text{med}\{P\}$ ,  $P$  is an unknown distribution over  $\{0, 1, \dots, 9\}$ .  $\vec{Y} = (4, 8, 9, 6, 2)$ .

## 6.1 Residual Bootstrap

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.**

### Example: Linear

Observation  $Y_i = a + b \frac{i}{n} + Z_i$ , where  $Z_i \sim N(0, \sigma^2)$  (i.i.d.) for  $i = 1, 2, \dots, n$

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

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

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

$$E_i = Y_i - \hat{a}_n - \hat{b}_n \frac{i}{n} \approx Z_i$$

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

$\Rightarrow$  Compute pseudo-data  $Y_i^{*(j)} = \hat{a}_n + \hat{b}_n \frac{i}{n} + E_i^{*(j)}$

$\Rightarrow$  Compute LS estimator

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

$\Rightarrow$  Evaluate bias

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

### Example: Nonlinear Markov Process

Observation  $Y_i = F_\theta(Y_{i-1}) + Z_i$ , where  $Z_i \sim N(0, \sigma^2)$  (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.)

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

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

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

$\Rightarrow$  Evaluate bias

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

# Chapter 7 Particle Filtering

Kalman filtering is used in tracking problems (dynamic models). Particle Filtering is an extension of Kalman filtering.

## 7.1 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.

## 7.2 Particle Filtering (Nonlinear Dynamic System)

Particle filtering is a nonlinear form of Kalman filtering, which doesn't have closed-from solutions.

We consider a Nonlinear Dynamic System

$$X_{t+1} \sim q(\cdot | X_t)$$

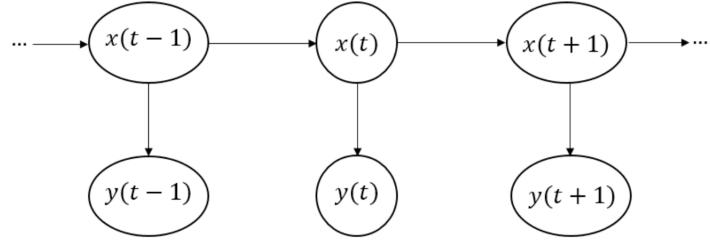
$$Y_t \sim r(\cdot | X_t)$$

$$t = 0, 1, 2, \dots$$

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 7.1:** Hidden Markov Model

### 7.2.1 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  $P(X_t|Y_{0:t-1})$

$$\begin{aligned}P(X_t|Y_{0:t}) &= P(X_t|Y_t, Y_{0:t-1}) \\ &= \frac{P(Y_t|X_t, Y_{0:t-1}) P(X_t|Y_{0:t-1})}{P(Y_t|Y_{0:t-1})} \\ &= \frac{r(Y_t|X_t) P(X_t|Y_{0:t-1})}{\int_{\mathbb{R}^m} r(Y_t|X_t) P(X_t|Y_{0:t-1}) dx_t}\end{aligned}$$

### 7.2.2 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 8 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  $\theta$  may be multidimensional, all of which are problematic for any numerical algorithm.

## Maximum-Likelihood (ML) Estimation

Given a vector  $\vec{y}$ , find the  $\theta$  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.

## 8.1 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.

$$P(y, z|\theta) = P(z|\theta), \quad \forall z \in h^{-1}(y)$$

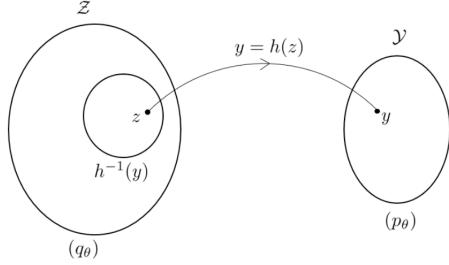
5. We can infer that the relationship between  $P(z|\theta)$  and  $P(y|\theta)$  is

$$P(z|\theta) = P(z|y, \theta)P(y|\theta), \quad \forall z \in h^{-1}(y)$$

$$P(y|\theta) = \sum_{z \in h^{-1}(y)} P(z|\theta), \quad \forall y$$

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 8.1:** 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  $\theta$ .

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  $\theta$  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 8.1

$\theta^*$  is a stable point of the EM algorithm if  $\exists$  subsequence that converges to  $\theta^*$ .

e.g.  $1, 3, \frac{1}{2}, 3, \frac{1}{3}, 3, \dots \frac{1}{n}, 3, \dots$



## 8.2 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  $\theta$ .

### 8.2.1 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  $\theta$  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,

$$\hat{\theta}_{ML} = \begin{cases} 0, & y^2 \leq 1 \\ y^2 - 1, & y^2 > 1 \end{cases}$$

### 8.2.2 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

$$\hat{\theta}^* = \begin{cases} 0, & y^2 \leq 1 \\ y^2 - 1, & y^2 > 1 \end{cases}$$

## 8.3 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  $\mu_j$ , and variances  $\sigma_j^2$ , 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  $\mu$  and variance  $\sigma^2$ .

### 8.3.1 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  $\theta$  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.

### 8.3.2 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  $\mu_j$ ,

$$\begin{aligned} 0 &= \frac{\partial Q(\theta|\hat{\theta}^{(k)})}{\partial \mu_j} = \sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j|y_i) \frac{(y_i - \mu_j)}{\sigma_j^2} \\ \hat{\mu}_j^{(k+1)} &= \frac{\sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j|y_i)y_i}{\sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j|y_i)} \end{aligned}$$

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)

### 8.3.3 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:**

$$\begin{aligned}\hat{\mu}_j^{(k+1)} &= \frac{\sum_{i=1}^n y_i \pi_{\hat{\theta}^{(k)}}(j | y_i)}{\sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j | y_i)} \\ (\hat{\sigma}_j^2)^{(k+1)} &= \frac{\sum_{i=1}^n (y_i - \hat{\mu}_j^{(k+1)})^2 \pi_{\hat{\theta}^{(k)}}(j | y_i)}{\sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j | y_i)} \\ \hat{\pi}^{(k+1)}(j) &= \frac{1}{n} \sum_{i=1}^n \pi_{\hat{\theta}^{(k)}}(j | y_i), \quad 1 \leq j \leq m.\end{aligned}$$

## 8.4 Convergence of EM Algorithm

### Theorem 8.1

The likelihood sequence  $p_{\hat{\theta}^{(k)}}(y)$ ,  $k = 0, 1, 2, \dots$  is nondecreasing.



### Proof 8.1

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 8.1

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  $\nabla \ln p_\theta = 0$ .



## 8.5 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:

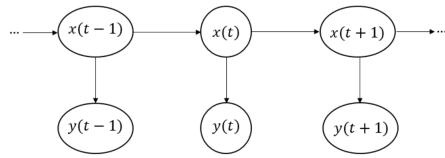
## Chapter 9 Hidden Markov model (HMM)

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  $\pi$ . 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 9.1:** 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)

### 9.1 Viterbi Algorithm: (MAP) estimate $X_{1:t}$ given $Y_{1:t}$

#### 9.1.1 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

$$\begin{aligned}\vec{x}^* &= \operatorname{argmax}_{\vec{x}} P(\vec{x}|\vec{y}) = \operatorname{argmax}_{\vec{x}} \frac{P(\vec{x}, \vec{y})}{P(\vec{y})} = \operatorname{argmax}_{\vec{x}} P(\vec{x}, \vec{y}) \\ &= \operatorname{argmax}_{\vec{x}} \pi(x_1) \prod_{t=1}^{n-1} A(x_t, x_{t+1}) \prod_{t=1}^n B(x_t, y_t) \\ &= \operatorname{argmax}_{\vec{x}} \ln \pi(x_1) + \sum_{t=1}^{n-1} \ln A(x_t, x_{t+1}) + \sum_{t=1}^n \ln B(x_t, y_t)\end{aligned}$$

### 9.1.2 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

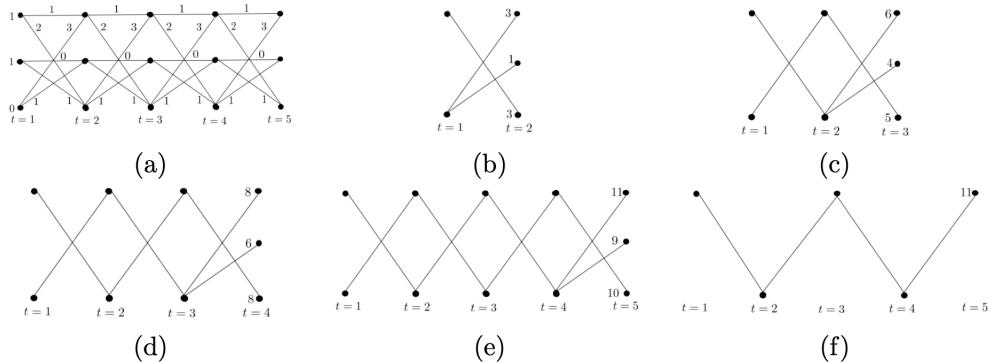
$$V(t, x_t = x) \triangleq \max_{x_1, x_2, \dots, x_{t-1}} \left[ \varepsilon([x_1, \dots, x_t]) = f(x_1) + \sum_{u=1}^{t-1} g_u(x_u, x_{u+1}) \right]$$

$$V(t, x_t = x) = \max_{x'} [V(t-1, x_{t-1} = x') + g_{t-1}(x', x)], t \geq 2$$

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 9.2:** (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  $\varepsilon(\vec{x}^*) = 11$ .

## 9.2 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

$$\hat{x}_t = \sum_{x_t \in X} x_t P(X_t = x | Y = \vec{y}), \quad 1 \leq t \leq n$$

## 9.3 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  $\gamma_t$  is the first marginal of  $\xi_t$ . The forward-backward algorithm allows efficient computation of these probabilities.

### 9.3.1 $\gamma_t(x) \triangleq P\left\{X_t = x \mid \vec{Y} = \vec{y}\right\}$

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

$$\mu_t(x) = P \left\{ Y_{1:t} = y_{1:t}, X_t = x \right\}, \quad x \in \mathcal{X}, 1 \leq t \leq n,$$

for which we derive a **forward recursion**. The recursion is initialized with

$$\mu_1(x) = P \left\{ Y_1 = y_1, X_1 = x \right\} = \pi(x)B(x, y_1).$$

For  $t \geq 1$  we express  $\mu_{t+1}$  in terms of  $\mu_t$  as follows:

$$\begin{aligned} \mu_{t+1}(x) &= P \left\{ Y_{1:t+1} = y_{1:t+1}, X_{t+1} = x \right\} \\ &\stackrel{(a)}{=} P \left\{ Y_{1:t} = y_{1:t}, X_{t+1} = x \right\} P \left\{ Y_{t+1} = y_{t+1} \mid X_{t+1} = x \right\} \\ &= B(x, y_{t+1}) \sum_{x' \in \mathcal{X}} P \left\{ Y_{1:t} = y_{1:t}, X_{t+1} = x, X_t = x' \right\} \\ &\stackrel{(b)}{=} B(x, y_{t+1}) \sum_{x' \in \mathcal{X}} P \left\{ Y_{1:t} = y_{1:t}, X_t = x' \right\} P \left\{ X_{t+1} = x \mid X_t = x' \right\} \\ &= B(x, y_{t+1}) \sum_{x' \in \mathcal{X}} \mu_t(x') A(x', x), \quad t = 1, 2, \dots, n-1 \end{aligned}$$

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

$$\nu_t(x) = P \left\{ Y_{t+1:n} = y_{t+1:n} \mid X_t = x \right\}, \quad x \in \mathcal{X}, 1 \leq t < n.$$

Starting from  $\nu_n(x) \equiv 1$ , we have the following **backward recursion**, expressing  $\nu_{t-1}$  in terms of  $\nu_t$  for  $2 \leq t \leq n$ :

$$\begin{aligned} \nu_{t-1}(x) &= P \left\{ Y_{t:n} = y_{t:n} \mid X_{t-1} = x \right\} \\ &= \sum_{x' \in \mathcal{X}} P \left\{ Y_{t:n} = y_{t:n}, X_t = x' \mid X_{t-1} = x \right\} \\ &\stackrel{(a)}{=} \sum_{x' \in \mathcal{X}} P \left\{ Y_{t:n} = y_{t:n} \mid X_t = x' \right\} P \left\{ X_t = x' \mid X_{t-1} = x \right\} \\ &\stackrel{(b)}{=} \sum_{x' \in \mathcal{X}} P \left\{ Y_{t+1:n} = y_{t+1:n} \mid X_t = x' \right\} P \left\{ Y_t = y_t \mid X_t = x' \right\} P \left\{ X_t = x' \mid X_{t-1} = x \right\} \\ &= \sum_{x' \in \mathcal{X}} \nu_t(x') B(x', y_t) A(x, x'), \quad t = n, n-1, \dots, 2 \end{aligned}$$

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.

$$9.3.2 \quad \xi_t(x, x') \triangleq P\left\{X_t = x, X_{t+1} = x' \mid \vec{Y} = \vec{y}\right\}$$

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

$$\xi_t(x, x') = \frac{\mu_t(x) A(x, x') B(x', y_{t+1}) \nu_{t+1}(x')}{\sum_{x, x' \in \mathcal{X}} \mu_t(x) A(x, x') B(x', y_{t+1}) \nu_{t+1}(x')}, \quad 1 \leq t \leq n, x, x' \in \mathcal{X}$$

### 9.3.3 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  $\alpha_t$  and  $c_t$ , and a backward recursion for  $\beta_t$ .

The time and storage complexity of the algorithm is  $O(n|\mathcal{X}|^2)$ .

# Chapter 10 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.

## 10.1 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.

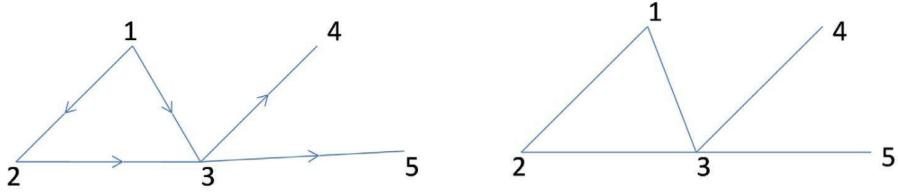
8. **Cycle/Loop:** An  $n$ -cycle, or loop, is a path of length  $n$   $i \rightarrow j$  with  $i = j$ .

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 10.1:** (a) Directed and (b) Undirected graph.

## 10.2 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.

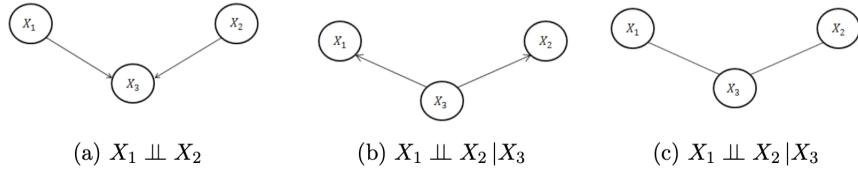
For instance a Markov chain is a chain-type directed acyclic graph where  $V = \{1, 2, \dots, n\}$ , and  $\pi(v) = v - 1$  for  $v \geq 2$ . The pmf for the sequence  $\vec{x}$  is obtained from the chain rule

$$p(\vec{x}) = p(x_1)p(x_2|x_1) \cdots p(x_n|x_{n-1})$$

## 10.3 Markov Networks

### 10.3.1 General Form

We can use undirected graph to represent conditionally independent.



**Figure 10.2:** (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  $\Omega$  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

### 10.3.2 Hammersley-Clifford theorem

#### Theorem 10.1 (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  $G$



### 10.3.3 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.

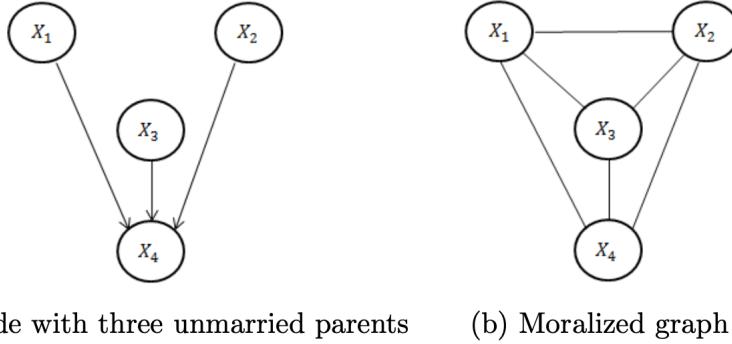
$$p(\vec{x}) = \frac{1}{Z} e^{-\sum_C V_C(x_C)}$$

This probability follows **Gibbs distribution (Boltzmann distribution)**. This distribution follows exponential families.

## 10.4 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.



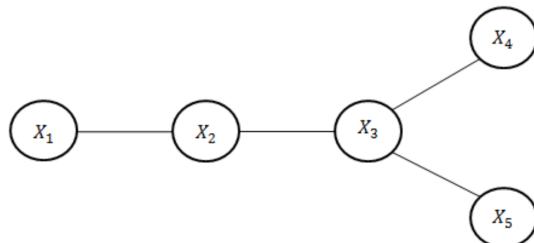
**Figure 10.3:** Graph moralization

## 10.5 Inference and Learning

### 10.5.1 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 10.4:** 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

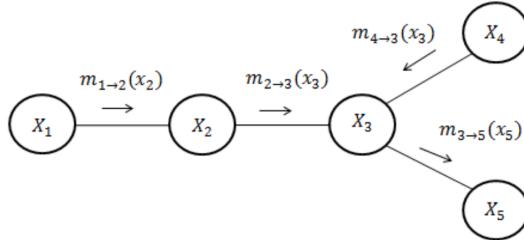
infeasible for large graphs (the number of items in the sum is  $|\mathcal{X}|^4$ );

$$p(x_5) = \sum_{x_1, \dots, x_4} \frac{1}{Z} \psi_{12}(x_1, x_2) \psi_{23}(x_2, x_3) \psi_{34}(x_3, x_4) \psi_{35}(x_3, x_5)$$

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 10.5:** 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  $p(x_2, x_5)$

$$p(x_2, x_5) = \frac{1}{Z} \sum_{x_3} \psi_{23}(x_2, x_3) \psi_{35}(x_3, x_5) \underbrace{\sum_{x_4} \psi_{34}(x_3, x_4)}_{m_{4 \rightarrow 3}(x_3)} \underbrace{\sum_{x_1} \psi_{12}(x_1, x_2)}_{m_{1 \rightarrow 2}(x_2)}.$$

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 11 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  $q$ .

## 11.1 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:

$$q_i^{(k+1)}(x_i) = \frac{1}{Z_i} \exp \left( \mathbb{E}_{\prod_{j \neq i}^n q_j^{(k)}} [\ln p(X_{1:i-1}, x_i, X_{i+1:n})] \right)$$

### 11.1.1 Graphical Models

Consider  $P$  = 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,

$$q_i(x_i) = \frac{1}{Z_i} \exp \left( \sum_{j \in N(i)} \sum_{x_j \in X} q_j(x_j) \ln \psi_{ij}(x_i, x_j) \right)$$

### 11.1.2 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  $\beta$  represents the inverse of a temperature. For  $\beta = 0$  the distribution is uniform, hence fully factorized. For large positive values of  $\beta$ , 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  $\beta$ , 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.

## 11.2 Exponential Families of Probability Distributions

### Definition 11.1 ( $d$ -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  $p_\theta$  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:**  $P_\theta = N(0, \theta^{-1})$ ,  $\theta$  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:

$$\Theta = \{\theta : \int_X e^{\theta^T T(y)} dy < \infty\}$$

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  $\mu$  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  $\mu$  is a correlation matrix, and so  $\mathcal{M}$  is the set of all  $n \times n$  symmetric nonnegative definite matrices.

## 11.3 ML Estimation

Consider  $n$  iid samples  $X^{(i)}, 1 \leq i \leq n$  drawn from the exponential distribution  $p_\theta$ . The ML estimator of  $\theta$  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.

## 11.4 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  $\lambda_{d+1}$  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  $\lambda$  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.

## 11.5

$$\begin{aligned} \min_{q \in Q} D(q||p) &= \min_q \mathbb{E}_q \left[ \ln \frac{q(x)}{p_\theta(x)} \right] \\ &= \min_q [A(\theta) - \theta^T \mathbb{E}_q[T(x)] - H(q)] \\ &= A(\theta) - \max_{\mu} \max_{q: \mathbb{E}_q[T(x)]=\mu} [\theta^T \mathbb{E}_q[T(x)] + H(q)] \\ &= A(\theta) - \max_{\mu \in M} \left[ \theta^T \mu + \max_{q: \mathbb{E}_q[T(x)]=\mu} H(q) \right] \end{aligned}$$

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).  $A^*(\mu) = \max_{\theta}[\theta^T \mu - A(\theta)]$ ,  $A(\theta) = \max_{\mu}[\theta^T \mu - A^*(\mu)]$

$$\min_{q \in Q} D(q \| p) = A(\theta) - \max_{\mu \in M} [\theta^T \mu + A^*(\mu)]$$

## 11.6 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  $d = |V||X| + |E||X|^2$ .

### 11.6.1 Marginal polytope

#### Definition 11.2 (marginal polytope)

*The mean parameters associated with the distribution  $p_\theta$  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  $\mathcal{M}(G)$ .*

### 11.6.2 Locally Consistent Marginal Distributions

Given a graph  $G = (V, E)$ , consider the set of marginal distributions  $\tau_i$  on individual nodes  $i \in V$  and pairwise marginals  $\tau_{ij}$  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 11.3 (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 11.1 (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 11.1**

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  $\mu$ , hence  $\mu \in \mathcal{M}(\mathcal{G})$ . This proves the first part of the claim.

**Example 11.1** 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  $\tau_1, \tau_2, \tau_3$  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  $\varepsilon$  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  $\tau$  may not be the marginals of any joint distribution on  $\mathcal{G}$ ,  $\tau$  are often referred to as **pseudomarginals**.

#### Definition 11.4 (pseudomarginal)

Marginal distribution  $\tau$  such that  $\tau \in \mathcal{L}(\mathcal{G})$  and  $\tau \notin \mathcal{M}(\mathcal{G})$  is called **pseudomarginals**.



### 11.6.3 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  $\mu_{ij}$ . 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  $\mu$ . Equivalently, the set of distributions  $p$  on a tree

graph is generally nonconvex. (Consider a length-3 chain for instance.)

#### 11.6.4 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

$$0 = \frac{\partial L(\mu, \lambda)}{\partial \mu_i(x)} \Rightarrow \mu_{i(x)} = \frac{1}{Z} e^{\theta_{ix} + \sum_{i \in N(i)} \sum_{x' \in X} \theta_{ijxx'} \mu_j(x')}$$

#### 11.6.5 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  $G' = (V, E')$ )

$$\mu_{12}(x_1, x_3) = \sum_{x_2} p(x_1, x_2, x_3) = \sum_{x_2} p(x_1, x_2)p(x_3|x_2) = \sum_{x_2} p(x_1, x_2) \frac{p_{23}(x_2, x_3)}{p(x_2)}$$

#### 11.6.6 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 11.5 (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 11.2 (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 12 $\ell_1$ Penalized Least Squares Minimization

We will focus on an  $\ell_1$ -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  $\ell_1$  norm of an unknown signal to be recovered. The first term is smooth, the second is not.

## 12.1 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  $\ell_2$  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  $\ell_1$  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  $\sigma^2$ ,

$$y = Ax + z, \text{ where } z \sim N(0, \sigma^2)$$

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.

## 12.2 Special Cases

### 12.2.1 Definition: Soft Threshold

#### Definition 12.1 (Soft Threshold)

$$S_\lambda(y) \triangleq \begin{cases} y - \lambda, & \text{for } y \geq \lambda \\ y + \lambda & \text{for } y \leq -\lambda \\ 0 & \text{for } |y| < \lambda \end{cases}$$



### 12.2.2 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,

$$0 = f'(x) \Rightarrow x = S_\lambda(y) \triangleq \begin{cases} y - \lambda, & \text{for } y \geq \lambda \\ y + \lambda & \text{for } y \leq -\lambda \\ 0 & \text{for } |y| < \lambda \end{cases}$$

### 12.2.3 Orthonormal A

If  $m = n$  and  $A$  is orthonormal, then  $A^{-1} = A^T$  and

$$\|y - Ax\|^2 = \|A(A^T y - x)\|^2 = \|A^T y - x\|^2$$

$$\Rightarrow x = S_\lambda(A^T y)$$

### 12.2.4 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  $\alpha$  is the step size.

## 12.3 General Solution: Lasso

## 12.4 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.

### 12.4.1 Proximal Minimization Algorithm

#### Definition 12.2 (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 12.1 (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  $F$ .



### 12.4.2 Apply to $\ell_1$ -penalized least-squares

To apply this strategy to  $\ell_1$ -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  $\alpha = \frac{1}{c}$ ).

## 12.5 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  $\alpha$  and  $\beta$  are optimized.

## 12.6 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.

## 12.7 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 13 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.

## 13.1 Definitions related to Sparsity

**Definitions:**

1.

### Definition 13.1 ( $\ell_p$ norm)

The  $\ell_p$  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 13.2 ( $\ell_0$ pseudonorm)

The  $\ell_0$  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 13.3 ( $k$ -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 13.4 (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  $\Sigma_1$  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 13.1 (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 13.5 ( $\ell_p$  approximation error)**

*The  $\ell_p$  approximation error of  $x \in \mathbb{R}^n$  in  $\Sigma_k$  is*

$$e_{k,p}(x) \triangleq \min_{\hat{x} \in \Sigma_k} \|x - \hat{x}\|_p, \quad p \geq 1$$



**Lemma 13.1 (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  $\Theta$  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  $\sigma_j$  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).

## 13.2 Measurement Matrix

### Definition 13.6 (Measurement Matrix)

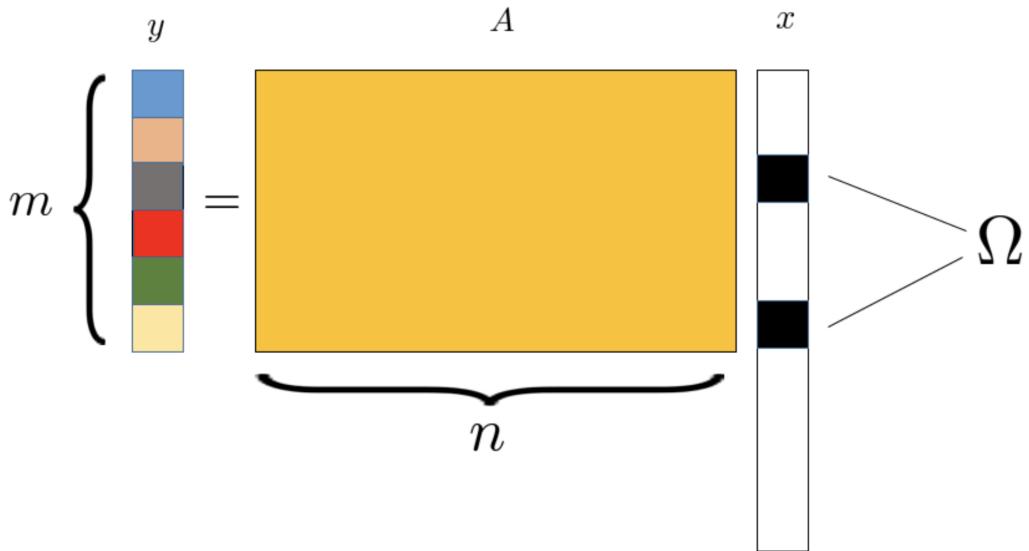
The measurement matrix  $A$  is a  $m \times n$  matrix where  $m \ll n$  (fat matrix). The observations are given by

$$\vec{y} = A\vec{x} \in \mathbb{R}^m \text{ in case of noise-free measurements.}$$



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 13.1:** Measurement of sparse signal  $\vec{x}$  with support set  $\Omega$  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

$$\vec{y} = A\vec{x} = A_\Omega\vec{x}_\Omega$$

### 13.2.1 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  $\sigma_k$ ,  $1 \leq k \leq m$ , 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 13.7 (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 13.8 (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.

### 13.2.2 Recovery of $k$ -Sparse Signals

#### Lemma 13.2

Unique recovery of  $x \in \Sigma_k$  given measurement matrix  $A$  is possible  $\Rightarrow \Sigma_{2k} \cap N(A) = \emptyset$ .



#### Proof 13.1

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 13.1

Unique recovery of  $\vec{x} \in \Sigma_k$  given measurement matrix  $A$  is possible  $\Leftrightarrow \text{spark}(A) > 2k$ . ( $q \geq 2k$ )



#### Proof 13.2

Eldar, Y. C., & Kutyniok, G. (Eds.). (2012). Compressed sensing: theory and applications, Chapter 1.

### 13.2.3 Restricted Isometry Property

#### Definition 13.9 (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  $\ell_2$  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 13.2 (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  $\Omega$  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 13.2 (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  $c = \frac{1}{2} \ln(1 + \sqrt{24}) \approx 0.28$ .



## 13.3 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  $\ell_0$  recovery problem

$$\min_{\vec{x} \in \mathbb{R}^n} \|\vec{x}\|_0 \quad \text{subj. to} \quad A\vec{x} = \vec{y}$$

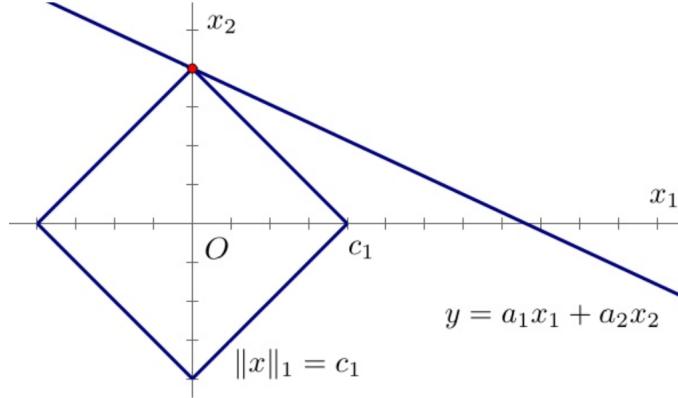
Unfortunately this problem is highly nonconvex. Solving it essentially requires evaluating all possible support sets of  $\vec{x}$ , which is a combinatorial problem.

A reasonable substitute is the so-called  $\ell_1$  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

$\ell_1$  ball  $\|\mathbf{x}\|_1 = \text{cst}$  at the solution  $\mathbf{x}$ , producing a 1-sparse solution.



**Figure 13.2:**  $\ell_1$  recovery for  $n = 2$  and  $m = 1$

The following fundamental theorem shows that the  $\ell_1$  recovery procedure is remarkably good.

**Theorem 13.3 ( $\ell_1$  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  $\ell_1$  approximation error of  $x$  in  $\Sigma_k$ .



**Proof 13.3**

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 13.1**

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.

## 13.4 Robust Signal Recovery from Noisy Observations

### 13.4.1 Bounded Noise

We consider observations corrupted by bounded noise:  $y = Ax + z$  where  $\|z\|_2 \leq \epsilon$ . We study the  $\ell_1$  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 13.4

*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).