

# Lecture 18

## Last time

- ▷ Matrix Calculus
- ▷ Matrix Chernoff
- ▷ Matrix Hoeffding

## Today

- ▷ Intro to estimation
- ▷ Maximum likelihood estimation

## Intro to estimation

So far we have developed foundational tools, in the time we have left we leverage these tools for two important statistical tasks: estimation and learning.

We start talking about estimation.

Suppose we are given a sample

$$z_1, \dots, z_n \sim P_\theta \leftarrow \text{Supported in } \mathbb{R}^d$$

where  $P_\theta$  is a distribution that depends on some object  $\theta$

**Goal:** Estimate  $\theta$  accurately.

**Examples:** Mean, Covariance, CDF,

There are two types of estimation:

▷ Parametric: Suppose that  $\theta \in \mathbb{R}^k$  and we can describe  $P_\theta$  via a finite-dimensional parameterization involving  $\theta$ .

Example: We wish to estimate  $\theta \in \mathbb{R}^d$  and assume

$$z_i \sim N(\theta, I).$$

▷ Nonparametric:  $P_\theta$  does not belong to some "nice" finite dimensional family.

Example: The object we want to estimate  $\theta$  is the density of the distribution and it can be any continuous function.

We will focus on parametric problems.

Def: An estimator of  $\theta$  is a (measurable) function of the data

$$\hat{\theta}_n = g(x_1, \dots, x_n).$$

We say that an estimator is unbiased if

$$E \hat{\theta}_n = \theta.$$

How do we measure accuracy?

One popular way is via the Mean Squared Error:

$$MSE(\hat{\theta}_n) := E \| \theta - \hat{\theta}_n \|_2^2.$$

In turn, this can be decomposed as

$$\begin{aligned} MSE(\hat{\theta}) &= E \| \theta - E \hat{\theta}_n + E \hat{\theta}_n - \hat{\theta}_n \|_2^2 \\ &= E \| \theta - E \hat{\theta}_n \|_2^2 \\ &\quad + E \langle \theta - E \hat{\theta}_n, E \hat{\theta}_n - \hat{\theta}_n \rangle \\ &\quad + E \| E \hat{\theta}_n - \hat{\theta} \|_2^2 \\ &= \underbrace{E \| \theta - E \hat{\theta}_n \|_2^2}_{\text{Bias}(\hat{\theta}_n)^2} + \underbrace{E \| \hat{\theta}_n - E \hat{\theta}_n \|_2^2}_{\text{Var}(\hat{\theta}_n)} \end{aligned}$$

Often, it pays to incur in some bias to reduce the overall MSE (see Stein's Paradox).

# Maximum Likelihood Estimation

MLE gives a systematic way to come up with estimators.

Suppose we have a parametric family

$$\{f(x; \theta) : \theta \in \Theta \subseteq \mathbb{R}^k\}.$$

Density for  $P_0$  (Probability of  $x$  for discrete dist)

Given that we observed  $z_1, \dots, z_n$ , the likelihood function is

$$L_n(\theta) = \prod_{i=1}^n f(z_i; \theta).$$

The log-likelihood is given by

$$\ln(\theta) = \sum_{i=1}^n \log f(z_i; \theta).$$

The MLE is given by

$$\hat{\theta}_n \in \underset{\theta \in \Theta}{\text{arg max}} \ln(\theta).$$

Intuitively, this corresponds to estimating the parameter that makes the data more likely.

Note that we can equivalently have

$$\hat{\Theta}_{\text{MSE}} \in \arg \max_{\theta \in \Theta} \ln(\theta).$$

Let's work out a couple of examples.

**Example:** Suppose  $Z_1, \dots, Z_n \sim \text{Ber}(p)$

Then

$$f(z; p) = p^z (1-p)^{1-z}.$$

Therefore

$$\begin{aligned} \ln(p) &= \sum_{i=1}^n z_i \log p + \sum_{i=1}^n (1 - z_i) \log(1-p) \\ &= S_n \log p + (n - S_n) \log(1-p) \end{aligned}$$

This function is concave in  $p$  and  
 $\nabla \ln(p) = 0$  reveals that

$$\hat{p}_{\text{MLE}} = \frac{1}{n} S_n.$$

-

**Example:** Suppose that  $\theta \in \mathbb{R}^d$  and we observe iid samples

$$Z_i = (X_i, Y_i) \text{ with}$$

$$Y_i = \theta^T X_i + \varepsilon_i. \quad \varepsilon_i \sim N(0, \sigma^2)$$

$X_i \sim D$  or deterministic

Equivalently,  
 $y_i \sim N(\bar{\theta}^T x_i, \sigma^2)$ .

Our goal is to estimate  $\bar{\theta}$ . Notice that there is another parameter that we don't care about, i.e., the variance  $\sigma^2$ . This is known as a nuisance parameter. Often (somewhat magically) the MSE doesn't require to know nuisance parameters, e.g.,

$$\begin{aligned} \ln(\theta) &= \frac{1}{(2\pi\sigma^2)^{n/2}} \prod_{i=1}^n e^{-\frac{(y_i - \theta^T x_i)^2}{2\sigma^2}} \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum (y_i - \theta^T x_i)^2\right) \end{aligned}$$

$$\Rightarrow \ln(\theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \theta^T x_i)^2$$

Notice that maximizing  $\ln(\theta)$  is equivalent to

*Why?*

$$\hat{\theta}_{MLE} = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{2n} \sum_{i=1}^n (y_i - \theta^T x_i)^2$$

which recovers least squares.

Beyond just estimating a parameter we might want to get confidence intervals or use the sample estimate for hypothesis testing; all of these are dubbed as "statistical inference" problems; take Statistical Theory I & II for more.

## MSE via Stochastic Optimization

Notice that we can always write

$$\hat{\theta}_{\text{MSE}} \in \operatorname{argmin}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log(f(z_i; \theta)).$$

Stochastic optimization aims to solve problems of the form

$$(P) \quad \min_{\theta \in \Theta} \mathbb{E}_{z \sim D} F(\theta, z)$$

$F: \mathbb{R}^k \times \mathbb{R}^m \rightarrow \mathbb{R}$

having access only to an iid sample  $z_1, \dots, z_n \sim D$ . Since we don't have access to  $D$ , a natural approach is to solve

$$(P_n) \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n F(\theta, z_i).$$

This is known as Empirical risk minimization. Clearly, MSE is an instance of ERM.

There are a number of interesting questions:

- ▷ How do we solve  $(P_n)$ ?  
Take Nonlinear I and II.
- ▷ Say  $\theta_n^*$  is a solution of  $(P_n)$  and  $\theta^*$  is a solution of  $(P)$ . Do we have that  $\theta_n^* \xrightarrow{\text{a.s.}} \theta^*$ ? How about convergence rates?  
We will study this question for linear least squares.
- ▷ How about confidence intervals?  
High prob. guarantees?  
This is a potential project!