

# Lecture Notes 7

## Parametric Point Estimation

### 1 Introduction

$X_1, \dots, X_n \sim p(x; \theta)$ . Want to estimate  $\theta = (\theta_1, \dots, \theta_k)$ . An *estimator*

$$\hat{\theta} = \hat{\theta}_n = w(X_1, \dots, X_n)$$

is a function of the data. Keep in mind that the parameter is a fixed, unknown constant. The estimator is a random variable.

For now, we will discuss three methods of constructing estimators:

1. The Method of Moments (MOM)
2. Maximum likelihood (MLE)
3. Bayesian estimators.

Later we will discuss some other methods. We will also discuss several methods for evaluating estimators including:

1. Bias and Variance
2. Mean squared error (MSE)
3. Minimax Theory
4. Large sample theory.

**Some Terminology.** Throughout these notes, we will use the following terminology:

1.  $\mathbb{E}_\theta(\hat{\theta}) = \int \cdots \int \hat{\theta}(x_1, \dots, x_n) p(x_1; \theta) \cdots p(x_n; \theta) dx_1 \cdots dx_n$ .
2. **Bias**:  $\mathbb{E}_\theta(\hat{\theta}) - \theta$ .
3. The distribution of  $\hat{\theta}_n$  is called its *sampling distribution*.
4. The standard deviation of  $\hat{\theta}_n$  is called the *standard error* denoted by  $se(\hat{\theta}_n)$ .
5.  $\hat{\theta}_n$  is *consistent* if  $\hat{\theta}_n \xrightarrow{P} \theta$ .
6. Later we will see that if  $\text{bias} \rightarrow 0$  and  $\text{Var}(\hat{\theta}_n) \rightarrow 0$  as  $n \rightarrow \infty$  then  $\hat{\theta}_n$  is **consistent**.
7. An estimator is *robust* if it is not strongly affected by perturbations in the data (more later).

## 2 The Method of Moments

Suppose that  $\theta = (\theta_1, \dots, \theta_k)$ . Define

$$\begin{aligned} m_1 &= \frac{1}{n} \sum_{i=1}^n X_i, & \mu_1(\theta) &= \mathbb{E}(X_i) \\ m_2 &= \frac{1}{n} \sum_{i=1}^n X_i^2, & \mu_2(\theta) &= \mathbb{E}(X_i^2) \\ & \vdots & \vdots & \\ m_k &= \frac{1}{n} \sum_{i=1}^n X_i^k, & \mu_k(\theta) &= \mathbb{E}(X_i^k). \end{aligned}$$

Let  $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$  solve:

$$m_j = \mu_j(\hat{\theta}), \quad j = 1, \dots, k.$$

In other words, we equate the first  $k$  sample moments with the first  $k$  theoretical moments. This defines  $k$  equations with  $k$  unknowns.

**Example 1**  $N(\beta, \sigma^2)$  with  $\theta = (\beta, \sigma^2)$ . Then  $\mu_1 = \beta$  and  $\mu_2 = \sigma^2 + \beta^2$ . Equate:

$$\frac{1}{n} \sum_{i=1}^n X_i = \hat{\beta}, \quad \frac{1}{n} \sum_{i=1}^n X_i^2 = \hat{\sigma}^2 + \hat{\beta}^2$$

to get

$$\hat{\beta} = \bar{X}, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

**Example 2** Suppose

$$X_1, \dots, X_n \sim \text{Binomial}(k, p)$$

where both  $k$  and  $p$  are unknown. We get

$$kp = \bar{X}_n, \quad \frac{1}{n} \sum_{i=1}^n X_i^2 = kp(1-p) + k^2p^2$$

giving

$$\hat{p} = \frac{\bar{X}}{k}, \quad \hat{k} = \frac{\bar{X}^2}{\bar{X} - \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}.$$

The method of moments was popular many years ago because it is often easy to compute. Lately, it has attracted attention again. For example, there is a large literature on estimating “mixtures of Gaussians” using the method of moments.

### 3 Maximum Likelihood

The most popular method for estimating parameters is maximum likelihood. The reason is that, under certain conditions, the maximum likelihood estimator is optimal. This result was established by Sir Ronald Fisher and Lucian LeCam. We'll discuss optimality later.

The maximum likelihood estimator (mle)  $\hat{\theta}$  is defined as the maximizer of

$$L(\theta) = p(X_1, \dots, X_n; \theta) \stackrel{iid}{=} \prod_i p(X_i; \theta).$$

This is the same as maximizing the log-likelihood

$$\ell(\theta) = \log L(\theta).$$

Often it suffices to solve

$$\frac{\partial \ell(\theta)}{\partial \theta_j} = 0, \quad j = 1, \dots, k.$$

**Example 3** *Binomial.*  $L(p) = \prod_i p^{X_i} (1-p)^{1-X_i} = p^S (1-p)^{n-S}$  where  $S = \sum_i X_i$ . So

$$\ell(p) = S \log p + (n-S) \log(1-p)$$

and  $\hat{p} = \bar{X}$ .

**Example 4**  $X_1, \dots, X_n \sim N(\mu, 1)$ .

$$L(\mu) \propto \prod_i e^{-(X_i - \mu)^2/2} \propto e^{-n(\bar{X} - \mu)^2}, \quad \ell(\mu) = -\frac{n}{2}(\bar{X} - \mu)^2$$

and  $\hat{\mu} = \bar{X}$ . For  $N(\mu, \sigma^2)$  we have

$$L(\mu, \sigma^2) \propto \prod_i \frac{1}{\sigma} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \right\}$$

and

$$\ell(\mu, \sigma^2) = -n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2.$$

Set

$$\frac{\partial \ell}{\partial \mu} = 0, \quad \frac{\partial \ell}{\partial \sigma^2} = 0$$

to get

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

**Example 5** Let  $X_1, \dots, X_n \sim \text{Uniform}(0, \theta)$ . Then

$$L(\theta) = \frac{1}{\theta^n} I(\theta > X_{(n)})$$

and so  $\hat{\theta} = X_{(n)}$ .

Suppose that  $\theta = (\eta, \xi)$ . The *profile likelihood* for  $\eta$  is defined by

$$L(\eta) = \sup_{\xi} L(\eta, \xi).$$

To find the mle of  $\eta$  we can proceed in two ways. We could find the overall mle  $\hat{\theta} = (\hat{\eta}, \hat{\xi})$ . The mle for  $\eta$  is just the first coordinate of  $(\hat{\eta}, \hat{\xi})$ . Alternatively, we could find the maximizer of the profile likelihood. These give the same answer. Do you see why?

The mle is *equivariant*, if  $\eta = g(\theta)$  then  $\hat{\eta} = g(\hat{\theta})$ . Suppose  $g$  is invertible so  $\eta = g(\theta)$  and  $\theta = g^{-1}(\eta)$ . Define  $L^*(\eta) = L(\theta)$  where  $\theta = g^{-1}(\eta)$ . So, for any  $\eta$ ,

$$L^*(\hat{\eta}) = L(\hat{\theta}) \geq L(\theta) = L^*(\eta)$$

and hence  $\hat{\eta} = g(\hat{\theta})$  maximizes  $L^*(\eta)$ . For non invertible functions this is still true if we define

$$L^*(\eta) = \sup_{\theta: \tau(\theta)=\eta} L(\theta).$$

(In other words, the profile likelihood.)

**Example 6** *Binomial*. The mle is  $\hat{p} = \bar{X}$ . Let  $\psi = \log(p/(1-p))$ . Then  $\hat{\psi} = \log(\hat{p}/(1-\hat{p}))$ .

## 4 Bayes Estimator

To define the Bayes estimator, we begin by *treating  $\theta$  as a random variable*. This point requires much discussion (which we will have later). For now, just tentatively think of  $\theta$  as random. We start with a *prior distribution*  $p(\theta)$  on  $\theta$ . Note that

$$p(x_1, \dots, x_n | \theta) p(\theta) = p(x_1, \dots, x_n, \theta).$$

Now compute the *posterior distribution* by Bayes' theorem:

$$p(\theta | x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n | \theta) p(\theta)}{p(x_1, \dots, x_n)}$$

where

$$p(x_1, \dots, x_n) = \int p(x_1, \dots, x_n | \theta) p(\theta) d\theta.$$

This can be written as

$$p(\theta | x_1, \dots, x_n) \propto L(\theta) p(\theta) = \text{Likelihood} \times \text{prior}.$$

Now compute a point estimator from the posterior. For example:

$$\hat{\theta} = \mathbb{E}(\theta | x_1, \dots, x_n) = \int \theta p(\theta | x_1, \dots, x_n) d\theta = \frac{\int \theta p(x_1, \dots, x_n | \theta) p(\theta) d\theta}{\int p(x_1, \dots, x_n | \theta) p(\theta) d\theta}.$$

**This approach is controversial.** We will discuss the controversy and the meaning of the prior later in the course. For now, we just think of this as a way to define an estimator.

**Example 7** Let  $X_1, \dots, X_n \sim \text{Bernoulli}(\theta)$ . Let the prior be  $\theta \sim \text{Beta}(\alpha, \beta)$ . Hence

$$p(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)}$$

and

$$\Gamma(\alpha) = \int_0^\infty t^{\alpha-1} e^{-t} dt.$$

Set  $Y = \sum_i X_i$ . Then

$$p(\theta | X) \propto \underbrace{\theta^Y 1 - \theta^{n-Y}}_{\text{likelihood}} \times \underbrace{\theta^{\alpha-1} 1 - \theta^{\beta-1}}_{\text{prior}} \propto \theta^{Y+\alpha-1} 1 - \theta^{n-Y+\beta-1}.$$

Therefore,  $\theta | X \sim \text{Beta}(Y + \alpha, n - Y + \beta)$ . The Bayes estimator is

$$\tilde{\theta} = \frac{Y + \alpha}{(Y + \alpha) + (n - Y + \beta)} = \frac{Y + \alpha}{\alpha + \beta + n} = (1 - \lambda) \hat{\theta}_{mle} + \lambda \bar{\theta}$$

where

$$\bar{\theta} = \frac{\alpha}{\alpha + \beta}, \quad \lambda = \frac{\alpha + \beta}{\alpha + \beta + n}.$$

This is an example of a **conjugate prior**.

**Example 8** Let  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$  with  $\sigma^2$  known. Let  $\mu \sim N(m, \tau^2)$ . Then

$$\mathbb{E}(\mu | X) = \frac{\tau^2}{\tau^2 + \frac{\sigma^2}{n}} \bar{X} + \frac{\frac{\sigma^2}{n}}{\tau^2 + \frac{\sigma^2}{n}} m$$

and

$$\text{Var}(\mu | X) = \frac{\sigma^2 \tau^2 / n}{\tau^2 + \frac{\sigma^2}{n}}.$$

## 5 MSE

Now we discuss the **evaluation of estimators**. The **mean squared error (MSE)** is

$$\mathbb{E}_\theta(\hat{\theta} - \theta)^2 = \int \cdots \int (\hat{\theta}(x_1, \dots, x_n) - \theta)^2 p(x_1; \theta) \cdots p(x_n; \theta) dx_1 \cdots dx_n.$$

The **bias** is

$$B = \mathbb{E}_\theta(\hat{\theta}) - \theta$$

and the **variance** is

$$V = \text{Var}_\theta(\hat{\theta}).$$

**Theorem 9** *We have*

$$MSE = B^2 + V.$$

**Proof.** Let  $m = \mathbb{E}_\theta(\hat{\theta})$ . Then

$$\begin{aligned} MSE &= \mathbb{E}_\theta(\hat{\theta} - \theta)^2 = \mathbb{E}_\theta(\hat{\theta} - m + m - \theta)^2 \\ &= \mathbb{E}_\theta(\hat{\theta} - m)^2 + (m - \theta)^2 + 2\mathbb{E}_\theta(\hat{\theta} - m)(m - \theta) \\ &= \mathbb{E}_\theta(\hat{\theta} - m)^2 + (m - \theta)^2 = V + B^2. \end{aligned}$$

■

An estimator is **unbiased** if the bias is 0. In that case, the  $MSE = \text{Variance}$ . There is **often a tradeoff between bias and variance**. So low bias can imply high variance and vice versa.

**Example 10** Let  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ . Then

$$\mathbb{E}(\bar{X}) = \mu, \quad \mathbb{E}(S^2) = \sigma^2.$$

The MSE's are

$$\mathbb{E}(\bar{X} - \mu)^2 = \frac{\sigma^2}{n}, \quad \mathbb{E}(S^2 - \sigma^2)^2 = \frac{2\sigma^4}{n-1}.$$

We would like to choose an estimator with small MSE. However, the MSE is a function of  $\theta$ . Later, we shall discuss **minimax estimators**, that use the maximum of the MSE over  $\theta$  as a way to compare estimators.

## 6 Best Unbiased Estimators

What is the smallest variance of an unbiased estimator? This was once considered an important question. **Today we consider it not so important.** **There is no reason to require an estimator to be unbiased.** Having **small MSE is more important**. However, for completeness, we will briefly consider the question.

An estimator  $W$  is **UMVUE (Uniform Minimum Variance Unbiased Estimator)** for  $\tau(\theta)$  if (i)  $E_\theta(W) = \tau(\theta)$  for all  $\theta$  and (ii) if  $E_\theta(W') = \tau(\theta)$  for all  $\theta$  then  $\text{Var}_\theta(W) \leq \text{Var}_\theta(W')$ .

**The Cramer-Rao inequality** gives a lower bound on the variance of any unbiased estimator. The bound is:

$$\text{Var}_\theta(W) \geq \frac{\left(\frac{d}{d\theta} E_\theta W\right)^2}{E_\theta \left(\left(\frac{\partial}{\partial \theta} \log f(X; \theta)\right)^2\right)} = \frac{(\tau'(\theta))^2}{I_n(\theta)}.$$

There is also a link with sufficiency.

**Theorem 11 The Rao-Blackwell Theorem.** *Let  $W$  be an unbiased estimator of  $\tau(\theta)$  and let  $T$  be a sufficient statistic. Define  $W' = \phi(T) = E(W|T)$ . Then  $W'$  is unbiased and  $\text{Var}_\theta(W') \leq \text{Var}_\theta(W)$  for all  $\theta$ .*

Note that  $\phi$  is a well-defined estimator since, by sufficiency, it does not depend on  $\theta$ .

**Proof.** We have

$$E_\theta(W') = E_\theta(E(W|T)) = E_\theta(W) = \tau(\theta)$$

so  $W'$  is unbiased. Also,

$$\begin{aligned} \text{Var}_\theta(W) &= \text{Var}_\theta(E(W|T)) + E_\theta(\text{Var}(W|T)) \\ &= \text{Var}_\theta(W') + E_\theta(\text{Var}(W|T)) \\ &\geq \text{Var}_\theta(W'). \end{aligned}$$

■

**NOTE:** Many books discuss certain topics at this point such as: completeness, ancillarity etc. This stuff is pretty useless. We won't cover these topics.