# Estimation (Chapter 7)

# **General Problem:**

Model:  $\{P_{\theta} : \theta \in \Theta\}$ 

Observe  $X \sim P_{\theta}, \theta \in \Theta, \theta$  unknown.

Estimate  $\theta$ . (Pick a plausible distn from family.)

Or estimate  $\tau = \tau(\theta)$ .

Example:  $\theta = (\mu, \sigma^2), \tau(\theta) = \mu - \sigma$ .

# Terminology

A **point estimator** is a statistic W(X) (a function of the data). We usually (but not always) require that  $W(X) \in \Theta$ .

A **parameter** is a function  $\tau(\theta)$ . ("A" parameter  $\tau(\theta)$  is a function of "the" parameter  $\theta$ .)

If the data X is a random sample  $(X_1, \ldots, X_n)$ , these definitions correspond to those in elementary statistics:

A **statistic** is a characteristic of the sample.

A parameter is a characteristic of the population.

Notation: Point estimators of parameters  $\theta$  or  $\tau = \tau(\theta)$  are often designated  $\hat{\theta} = \hat{\theta}(X)$  or  $\hat{\tau} = \hat{\tau}(X)$ .

# **Examples of Parameters:**

Notation:

 $X = (X_1, \ldots, X_n)$  iid from the pdf (or pmf)  $f(x | \theta)$ .

X is a single rv from  $f(x \mid \theta)$ .

For concreteness, think of  $\theta = (\mu, \sigma^2)$  and  $X \sim N(\mu, \sigma^2)$ .

Some parameters:

$$\begin{aligned} \tau(\theta) &= \theta \\ \tau(\theta) &= \mu \text{ or } \tau(\theta) = \mu^2 \\ \tau(\theta) &= \sigma^2 \text{ or } \tau(\theta) = \sigma^4 \\ \tau(\theta) &= P_{\theta}(X \in A) = \int_A f(x \mid \theta) \\ \tau(\theta) &= E_{\theta}X = \int xf(x \mid \theta) \, dx \quad \text{(general case)} \\ \tau(\theta) &= E_{\theta}h(X) = \int h(x)f(x \mid \theta) \, dx \\ \tau(\theta) &= \text{median of } f(x \mid \theta) \\ \tau(\theta) &= \text{interquartile range of } f(x \mid \theta) \\ \tau(\theta) &= 95^{\text{th}} \text{ percentile of } f(x \mid \theta) \end{aligned}$$

# **Empirical Estimators:**

It is often possible to estimate a population quantity by a natural sample analog.

### **Examples:**

Parameter  $\tau(\theta)$ 

Estimate  $\hat{\tau}(X)$ 



†: F is the cdf of  $f(x | \theta)$  (the population cdf) and  $\hat{F}_n$  is the empirical cdf (defined later).

Intuitive approaches to estimation

Empirical Estimates (summary) Estimate a population quantity by the natural sample analog. For example, estimate population mean by sample mean, pop. variance by sample variance, pop. quantile by sample quantile, etc. Substitution Principle (Plug-in Method) Suppose  $\alpha = \alpha(\theta)$  and  $\beta = \beta(\theta)$  are two parameters related by  $\alpha = h(\beta)$ . If  $\hat{\beta} = \hat{\beta}(\chi)$  is a "reasonable" estimator of  $\beta$ , then  $\hat{\alpha} = h(\hat{\beta})$  is a "reasonable" estimator of  $\alpha$ . More generally, if a, B1, B2, ..., BK are parameters related by  $\alpha = h(\beta_1, \beta_2, \dots, \beta_K)$ , and B1, B2,..., BK are "reasonable" estimators of BI,...,BK, then  $\hat{\alpha} = h(\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k)$ 

is a "reasonable" estimator of  $\alpha$ .

Example: X1, X2,..., Xn iid N(11,02) X Estimate  $T(\theta) = P(-|\langle X < |)$ where  $X \sim N(\mu, \sigma^2)$ . (A) An empirical estimate  $\hat{\gamma} = \frac{1}{n} \sum_{i=1}^{n} I(-1 < X_i < 1)$ = the sample proportion  $Z = \frac{X - M}{T}$ (B) A Plug-in estimate  $= \oint (+1-\mu) - \oint (-1-\mu) \quad \text{where } \oint is$ cdf of N(O31)  $=h(\mu,\sigma)$ Reasonable estimates of M, or are  $\hat{\mu} = \overline{X}$ ,  $\hat{\sigma} = 5 = \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (X_i - \overline{X})^2}$ so that a plug-in estimate is given by  $\hat{\tau} = h(\hat{\mu}, \hat{\sigma}) = \Phi(\frac{1-\bar{x}}{\bar{\sigma}}) - \Phi(-\frac{1-\bar{x}}{\bar{\sigma}}).$ 

Which estimator is better? (A) or (B)? Intuition suggests (B) is better. Estimator (A) does not even use the assumption of normality. However, if it turns out that the normality assumption is false, then (A) may end up giving the better estimate of P(-I<X<I). Example: X1, X2,..., Xn iid from a Cauchy location-scale family with  $pdf \quad f(x|\mu,\sigma) = \frac{1}{\sigma} \cdot \frac{1}{\pi(1+px-\mu)^2}$ for  $-\infty < x < \infty$ . Estimate  $\Theta = (\mu, \sigma)$ . Note: This distribution does not have a finite mean. Thus  $\bar{x}$  and  $s^2$  are not useful here. Useful Facts: (where X~ f(·140))  $P(X < \mu) = .5$ Р(X< μ-σ)=.25  $P(X < \mu + \sigma) = .75$ 

Notation For 
$$0 ,
Let  $\beta_p = population p^{th} guantile$ ,  
 $Q_p = sample p^{th} guantile$ .  
(Formal definitions:  
Let  $F = pop. cdf$ :  $F(t) = P(x \le t)$   
 $\hat{F} = sample cdf$ :  $\hat{F}(t) = \frac{1}{n} \sum_{i=1}^{n} I(X_i \le t)$   
(empirical  
 $cdf$ )  
Then  
 $\beta_p = \inf\{x: F(x) \ge p\}$   
 $Q_p = \inf\{x: \hat{F}(x) \ge p\}$ .  
A "reasonable" estimate of  $\beta_p$  is  $\hat{\beta}_p = Q_p$ .  
The "useful facts" say that for the  
Cauchy L-S family  
 $\beta_{.5} = \mu$   
 $\beta_{.25} = \mu = \sigma$   
 $\beta_{.75} = \mu + \sigma$ )  
Thus  
 $\Theta = (\mu_s \sigma) = h(\beta_{.5}, \beta_{.25} - \beta_{.75})$   
 $= (\beta_{.5}, \frac{1}{2}(\beta_{.75} - \beta_{.25}))$   
so that a Plug-in estimate is given by  
 $\hat{\Theta} = h(\hat{\beta}_{.5}, \hat{\beta}_{.25}, \hat{\beta}_{.75}) = (Q_{.5}, \frac{1}{2}(Q_{.75} - Q_{.25})).$$$

Estimation by the Method of Moments (MOM) (Fitting distns by matching moments) MOM is a special case of the Plug-in method. Notation: (delete primes used in text)  $\mu_r = E X^r = r^{th} population moment$  $(M_r = M_r(\Theta) \text{ is a parameter.})$  $m_r = \frac{1}{n} \sum_{i=1}^{n} x_i^r = r^{th} \text{ sample moment}$ A "reasonable" estimate of  $\mu_r$  is  $\hat{\mu}_r = m_r$ . Thus .. MOM: If parameter  $\tau = \tau(\theta)$  can be written as a function of pop. moments  $\gamma = h(\mu_1, \mu_2, \dots, \mu_K),$ then a "reasonable" estimate of r is  $\hat{\tau} = h(m_1, m_2, \dots, m_K).$ 

# Parameter estimation by the Method of Moments Situation:

Suppose we have a model

 $X_1, X_2, \ldots, X_n$  iid  $f(x|\theta)$ 

where  $f(x|\theta)$  is the pdf (or pmf) of a family of distributions depending on a single parameter  $\theta$ .

The value of  $\theta$  is unknown.

We observe data  $x_1, x_2, \ldots, x_n$ .

How do we estimate  $\theta$ ?

### Notation:

Let X denote a single observation from  $f(x|\theta)$ .

Define

$$\mu = population mean = EX$$

$$\hat{\mu} = \bar{x} = \text{sample mean} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

Note that  $\mu$  is a function of  $\theta$ , say  $\mu = g(\theta)$ .

### Method of Moments (MOM):

Estimate  $\theta$  by that value  $\hat{\theta}$  which makes the population mean  $\mu$  equal to the sample mean  $\bar{x}$ .

#### Formal Procedure:

Step 1: Find  $\mu$  as a function of  $\theta$ :

 $\mu = EX = h(\theta).$ 

This is done either by looking up the family of distributions in the appendix or by doing the calculation

$$EX = \int_{-\infty}^{\infty} x f(x|\theta) dx$$
 or  $\sum_{\text{all } x} x f(x|\theta)$ .

Step 2: Solve for  $\theta$  as a function of  $\mu$ :

$$\theta = g(\mu) \,. \qquad (\dagger)$$

Step 3: Now plug in  $\hat{\mu} = \bar{x}$  to obtain the MOM estimate :  $\hat{\theta} = g(\bar{x})$ .

- Note: If  $\mu$  does not depend on  $\theta$  (for instance, if  $\mu = 0$  for all  $\theta$ ), then MOM is carried out using the second moment.
- Rationale: MOM works because the LLN guarantees that the sample mean  $\bar{x}$  will be close to the population mean  $\mu$  (with high probability) when the sample size n is large.

Since g (in  $\dagger$ ) is a continuous function,  $\bar{x} \approx \mu$  implies  $g(\bar{x}) \approx g(\mu)$  which says that  $\hat{\theta} \approx \theta$ .

Example: Suppose you observe  

$$x_{1}, x_{2}, ..., x_{n}$$
 iid Geometric (p).  
Find the MOM estimate of  $p$ .  
 $\square$  Find  $\mu$  as a function of  $p$ .  
 $\mu = EX = \sum_{\substack{x=1 \\ x=1}}^{\infty} x \cdot p(1-p)^{x-1} = \frac{1}{p}$   
not needed in  $f$   
from appendix  
 $\mu = \frac{1}{p}$   
(2) Solve for  $p$  as a function of  $\mu$ .  
 $p = \frac{1}{\mu}$   
(3) Plug in  $\hat{\mu} = \overline{x}$  for  $\mu$   
 $\hat{p} = \frac{1}{\lambda} = \frac{1}{\overline{x}}$   
Conclusion : The MOM estimate of  $p$  is  $\frac{1}{\overline{x}}$ .  
 $(\hat{p} = \frac{1}{\overline{x}})$ 

Parameter estimation by the Method of Moments

Situation: (multi-parameter case) Suppose we have a model

 $X_1, X_2, \ldots, X_n$  iid  $f(x|\theta)$ 

where  $f(x|\theta)$  is the pdf (or pmf) of a family of distributions depending on a vector of parameters

 $\theta = (\theta_1, \theta_2, \ldots, \theta_p).$ 

The vector of values  $\theta$  is unknown.

We observe data  $x_1, x_2, \ldots, x_n$ .

How do we estimate  $\theta$ ?

Notation:

Let X denote a single observation from  $f(x|\theta)$ .

Define

$$\mu_k = (population k-th moment) = EX^k$$

$$\widehat{\mu}_k = (\text{sample } k\text{-th moment}) = rac{1}{n} \sum_{i=1}^n x_i^k$$

(Special case:  $\mu = \mu_1$  and  $\bar{x} = \hat{\mu} = \hat{\mu}_1$ .)

Note that  $\mu_k$  is a function of  $\theta$ , say  $\mu_k = h_k(\theta_1, \ldots, \theta_p)$ .

### Method of Moments (MOM):

Estimate  $\theta = (\theta_1, \theta_2, \dots, \theta_p)$ by those values  $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_p)$ which make the population moments  $(\mu_1, \mu_2, \dots, \mu_p)$ equal to the sample moments  $(\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_p)$ .

MOM for  $\Theta = (\Theta_1, \Theta_2, \dots, \Theta_p)$ 1) Find expressions for u, Ma, ..., Mp:  $\mathcal{M}_{1} = h_{1}(\Theta_{1}, \ldots, \Theta_{p})$  $\mathcal{M}_2 = h_2(\Theta_1, \dots, \Theta_p)$  $\mu_{p} = h_{p}(\Theta_{1}, \dots, \Theta_{p})$ Look them up in appendix or evaluate using  $\mu_{k} = EX^{k} = \int x^{k} f(x|\theta) dx$  (continuous) or  $\sum_{\alpha \parallel x} \chi^{\kappa} f(x|\theta)$  (discrete) Solve this system of p equations for  $\Theta_1, \Theta_2, \dots, \Theta_p$ : 2)  $\Theta_1 = g_1(\mu_1, \dots, \mu_p)$  $\Theta_2 = g_2(\mu_1, \dots, \mu_p)$  $\Theta_{n} = g_{1p}(\mu_{1}, \dots, \mu_{p})$ 

3 Plug in  $\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_p$  as estimates of u,,..., up:  $\hat{\Theta}_1 = q_1(\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_p)$  $\hat{\theta}_2 = g_2(\hat{\mu}_1, \dots, \hat{\mu}_p)$  $\hat{\Theta}_{p} = g_{p}(\hat{\mu}_{1}, \dots, \hat{\mu}_{p}).$ Special Case: p=2 () Find My M2.  $\mu_1 = h_1(\Theta_1, \Theta_2)$  $\mathcal{M}_2 = h_2(\Theta_1, \Theta_2)$ (2) Solve for  $\Theta_{1}, \Theta_{2}$ :  $\Theta_1 = \mathcal{G}_1(\mathcal{U}_1, \mathcal{H}_2)$  $\Theta_2 = g_2(\mu_1,\mu_2)$ 3 Plug in  $\hat{\mu}_1, \hat{\mu}_2$  for  $\mu_1, \mu_2$ :  $\hat{\theta}_{1} = g_{1}(\hat{\mu}_{1}, \hat{\mu}_{2})$  $\hat{\Theta}_2 = g_2(\hat{\mu}_1, \hat{\mu}_2)$ 

Consistent Estimators A sequence of estimators  $W_h = W_h(X_1, X_2, ..., X_n)$ is a <u>consistent</u> sequence of estimators for the parameter  $T = T(\Theta)$  if, for every  $\varepsilon > 0$  and every  $\Theta \in \Theta$ ,  $\lim_{n \to \infty} \Theta(|W_n - T| < \varepsilon) = 1$ .

The sequence is <u>strongly consistent</u> if we may replace (\*) by  $W_n \rightarrow T$  with probability 1. (as  $n \rightarrow \infty$ ) The sequence is <u>consistent in 2<sup>nd</sup> mean</u> (or in L<sup>2</sup>) if we may replace (\*) by  $\lim_{n \rightarrow \infty} E_{\Theta} (W_n - T)^2 = 0$ . Let  $X_{1,X_{2}}, X_{3}, \dots$  be i.d. Strong Law of Large Numbers If  $E[h(X_{i})] < \infty$ , then  $\frac{1}{n} \sum_{i=1}^{n} h(X_{i}) \xrightarrow{\text{wp1}} Eh(X_{i})$  $as n \to \infty$ .



Another Fact: Suppose the population  $p^{\text{th}}$  quantile  $\beta_p$ is unique (that is, there exists a <u>unique</u> value  $\chi (=\beta_p)$  such that  $F(\chi) = p$  where F is the pop. cdf), then  $Q_p \xrightarrow{wp_1} \beta_p$ . (Sections 5.5 and 10.1 discuss modes of convergence and consistency of estimates in greater detail.)

The three types of consistency are (special cases of) 'convergence in probability', 'convergence almost surely', and 'convergence in  $L^2$ ' (or in mean square), respectively.

Preservation of convergence by continuous functions:

If  $W_n \to \tau$  in probability, and g is a continuous function, then  $g(W_n) \to g(\tau)$  in probability.

Also true for functions of many variables:

If  $U_n \to \xi$  and  $W_n \to \tau$  in probability, and  $g : \mathbb{R}^2 \to \mathbb{R}$  is a continuous function, then  $g(U_n, W_n) \to g(\xi, \tau)$  in probability.

The previous facts remain true if "in probability" is everywhere replaced by "almost surely".

Thus

Continuous functions of consistent estimates are consistent.

As a consequence, it is typically true that:

Estimators obtained by plug-in (substitution) are consistent.

Method of Moments (MOM) estimators are consistent.

Method of moments (MOM) estimators are (typically) continuous functions of sample moments, which are consistent estimates of populations moments.

# Example:

If  $X_1, X_2, X_3, \ldots$  are iid Geometric(*p*), the MOM estimator of *p* based on  $X_1, \ldots, X_n$  is  $1/\bar{x}_n$  where  $\bar{x}_n = n^{-1} \sum_{i=1}^n X_i$ .

WLLN implies  $\bar{x}_n \to EX_1 = 1/p$  in probability (as  $n \to \infty$ ).

Thus  $1/\bar{x}_n \rightarrow 1/(1/p) = p$  in probability.

This holds for all p. Thus  $1/\bar{x}_n$  is a consistent estimator of p.

Using the SLLN, the earlier statements remain true with 'in probability' replaced by 'almost surely' so that  $1/\bar{x}_n$  is also a *strongly* consistent estimator of p.

What about consistency in  $2^{nd}$  mean (or in  $L^2$ )?

**Example** Let  $X_1, X_2, X_3, \ldots$  are iid  $N(\mu, \sigma^2)$ . The most commonly used estimate of  $\sigma^2$  based on  $X_1, \ldots, X_n$  is

$$s_n^2 = (n-1)^{-1} \sum_{i=1}^n (X_i - \bar{x}_n)^2$$

 $s_n^2 \to \sigma^2$  in probability (for all  $\mu$  and  $\sigma^2$ ). (†)

Thus, applying the continuous function  $g(x) = \sqrt{x}$  to both sides:  $s_n \to \sigma$  in probability (for all  $\mu$  and  $\sigma^2$ ).

(These results don't require normality, but hold for any population with a finite second moment.)

Proof of  $(\dagger)$ :

Show that  $E(s_n^2 - \sigma^2)^2 = \operatorname{Var}(s_n^2) \to 0$ .

Alternatively, apply LLN to the identity:

$$s_n^2 = (n-1)^{-1} \left( \sum_{i=1}^n X_i^2 - n\bar{x}_n^2 \right) = \frac{n}{n-1} \left( \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{x}_n^2 \right)$$

$$\frac{E \times ample}{\text{with } pdf f(x)} = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha,\beta)}, \begin{array}{l} 0 < x < 1\\ (\alpha,\beta > 0) \end{array}$$
with  $pdf f(x) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha,\beta)}, \begin{array}{l} 0 < x < 1\\ (\alpha,\beta > 0) \end{array}$ 
where  $B(\alpha,\beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$ .
$$2 \text{ params } \alpha,\beta \Longrightarrow \text{ Use two moments.}$$

$$E \times = \mu_1 = \frac{\alpha}{\alpha+\beta}$$

$$E \times = \mu_1 = \frac{\alpha}{\alpha+\beta}$$

$$E \times = \mu_2 = \frac{\alpha(\alpha+1)}{(\alpha+\beta)(\alpha+\beta+1)} \quad \text{J} \text{ solve for } \alpha,\beta \text{ in } terms \text{ of } \beta,\beta \text{ in } \beta,\beta \text{ or } \beta,\beta \text{ in } \beta,\beta \text{ or }$$

Note that  $\mu_1 = \frac{\alpha}{\alpha + \beta} = \alpha d'$  $\Rightarrow \alpha = \mu_1 / \delta$  $\beta = (\alpha + \beta) - \alpha = \frac{1}{\alpha} - \frac{\mu_1}{\beta}$  $=\left|\frac{1}{\alpha}(1-\mu_{I})\right|=\beta$ 



 $= \frac{\mu_1 - \mu_2}{\mu_2 - \mu_1^2}$ 

In summary:

 $\alpha = \mu_1 \mathfrak{F}, \ \beta = (1 - \mu_1) \mathfrak{F}$ where  $S = \frac{1}{5} = \frac{\mu_1 - \mu_2}{\mu_1 - \mu_1^2}$ , so the MOM estimates are  $\hat{\alpha} = m_1 \hat{s}, \hat{\beta} = (1-m_1)\hat{s}, \hat{s} \equiv \frac{m_1 - m_2}{m_2 - m_1^2}$ where