

# 8. Estimation

## 8.1. Point estimation

- Point estimation is the use of the value  $\hat{\theta}$  of a statistic  $\hat{\theta} = g(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  to make a guess for the unknown parameter vector  $\theta \in \Theta \subset \mathbb{R}^K$  characterizing the distribution  $P_{\tilde{X}}(\cdot; \theta)$  of the random vector (or sample)  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$ .
- Note that, since the parameter  $\theta$  is unknown, the estimator  $\hat{\theta}$  is a statistic that cannot depend on  $\theta$  but only on the sample  $(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$ .
- The statistic  $\hat{\theta} = g(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  we use for estimation is called **estimator**. The value  $\hat{\theta} = g(x_1, x_2, \dots, x_n)$  taken by an estimator for a particular sample value  $X = (x_1, x_2, \dots, x_n)$  is called an **estimate**.
- For instance,  $\mathbf{s}^2 \neq \sigma^2$  in general. Questions: what is the relationship between  $\mathbf{s}^2$  and  $\sigma^2$ ?, how "close" are they?, does  $\mathbf{s}_n^2$  converge to  $\sigma^2$ ?, etc.

## 8.2. The mean square error of an estimator and the relative efficiency of estimators

- **Definition.** The mean square error (MSE) of an estimator  $\hat{\theta}$  for the population parameter  $\theta \in \mathbb{R}$  is

$$\mathbb{E} \left[ \left( \hat{\theta} - \theta \right)^2 \right].$$

- If  $\theta \in \mathbb{R}^K$ , then the MSE is

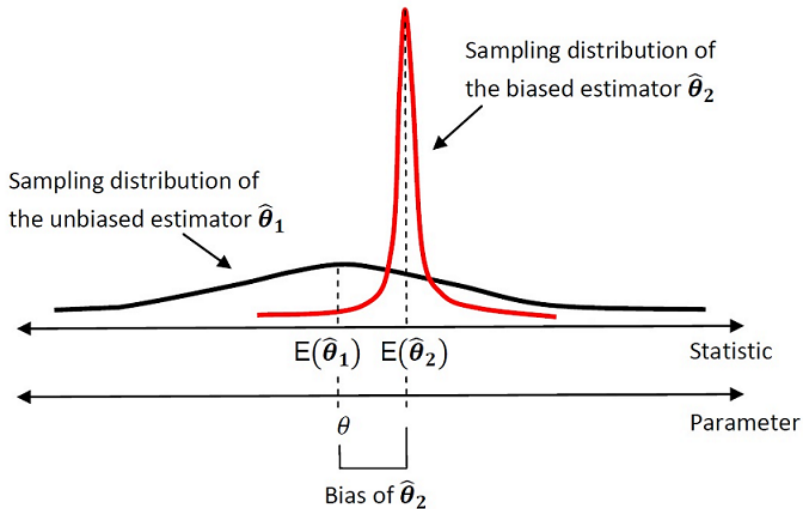
$$\mathbb{E} \left[ \left\| \hat{\theta} - \theta \right\|^2 \right].$$

- Note that

$$\mathbb{E} \left[ \left( \hat{\theta} - \theta \right)^2 \right] = \text{Var} \left( \hat{\theta} - \theta \right) + \left[ \mathbb{E} \left( \hat{\theta} - \theta \right) \right]^2 = \text{Var} \left( \hat{\theta} \right) + \left[ b_{\hat{\theta}} \left( \theta \right) \right]^2,$$

where  $b_{\hat{\theta}} \left( \theta \right) = \mathbb{E} \left( \hat{\theta} - \theta \right) = \mathbb{E} \left( \hat{\theta} \right) - \theta$  is the bias of the estimator  $\hat{\theta}$ .

A biased estimator with small variance may be preferable to an unbiased estimator with large variance



- **Definition.** We say that  $\hat{\theta}_1$  is a better (or a more efficient) estimator for  $\theta \in \mathbb{R}$  than  $\hat{\theta}_2$  if

$$\mathbb{E} \left[ \left( \hat{\theta}_1 - \theta \right)^2 \right] \leq \mathbb{E} \left[ \left( \hat{\theta}_2 - \theta \right)^2 \right].$$

- If the two previous estimators for  $\theta$  were unbiased,  $b_{\hat{\theta}_1}(\theta) = b_{\hat{\theta}_2}(\theta) = 0$ , then  $\hat{\theta}_1$  is better (or more efficient) than  $\hat{\theta}_2$  if

$$\text{Var} \left( \hat{\theta}_1 \right) \leq \text{Var} \left( \hat{\theta}_2 \right).$$

- **Definition.** We say that  $\hat{\theta}^*$  is the best (or the most efficient) estimator in the class  $C$  of estimators for  $\theta \in \mathbb{R}$  if

$$\mathbb{E} \left[ \left( \hat{\theta}^* - \theta \right)^2 \right] \leq \mathbb{E} \left[ \left( \hat{\theta} - \theta \right)^2 \right], \quad \text{for all } \hat{\theta} \in C.$$

## 8.3. The sample mean and sample variance as unbiased estimators

- Consider a random sample  $\{\tilde{x}_i\}_{i=1}^n$  of size  $n$  from a population  $\tilde{x}$  with distribution  $P_{\tilde{x}}$  having the mean  $\mu$  and the finite variance  $\sigma^2$ .
- Sample mean:

$$\bar{x} = \bar{\mathbf{x}}_n = \frac{\sum_{i=1}^n \tilde{x}_i}{n}.$$

- Sample variance:

$$\mathbf{s}^2 = \mathbf{s}_n^2 = \frac{\sum_{i=1}^n (\tilde{x}_i - \bar{\mathbf{x}}_n)^2}{n - 1}.$$

- Reminder:

**(a)**

$$E(\bar{\mathbf{x}}) = \mu \quad \text{and} \quad \text{Var}(\bar{\mathbf{x}}) = \frac{\sigma^2}{n}.$$

**(b)** Strong law of large numbers:

$$\bar{\mathbf{x}}_n \xrightarrow{a.s.} \mu.$$

**(c)** Central limit theorem:

$$\tilde{z}_n \equiv \frac{\bar{\mathbf{x}}_n - E(\bar{\mathbf{x}}_n)}{\sqrt{\text{Var}(\bar{\mathbf{x}}_n)}} = \frac{\bar{\mathbf{x}}_n - \mu}{\sigma / \sqrt{n}} \stackrel{a}{\sim} N(0, 1) \quad (\text{or } \tilde{z}_n \longrightarrow N(0, 1))$$

or, equivalently,

$$\sqrt{n}(\bar{\mathbf{x}}_n - \mu) \longrightarrow N(0, \sigma^2).$$

**(d)** If the population is normal, then

$$\bar{\mathbf{x}}_n \sim \mathbf{N}\left(\mu, \frac{\sigma^2}{n}\right) \iff \frac{\bar{\mathbf{x}}_n - \mu}{\sigma/\sqrt{n}} \sim \mathbf{N}(0, 1) \iff \sqrt{n}(\bar{\mathbf{x}}_n - \mu) \sim \mathbf{N}(0, \sigma^2).$$

**(e)**  $\mathbf{s}^2$  is an unbiased estimator for  $\sigma^2$  :

$$\mathbb{E}(\mathbf{s}^2) = \sigma^2.$$

**(f)** If the population mean  $\mu$  is known, then the statistic

$$\check{\mathbf{s}}^2 = \frac{\sum_{i=1}^n (\tilde{x}_i - \mu)^2}{n}$$

is an unbiased estimator for  $\sigma^2$ ,  $\mathbb{E}(\check{\mathbf{s}}^2) = \mathbb{E}(\mathbf{s}^2) = \sigma^2$ .

**(g)** If  $\hat{\mathbf{s}}^2 = \frac{\sum_{i=1}^n (\tilde{x}_i - \bar{\mathbf{x}})^2}{n} = \frac{(n-1)\mathbf{s}^2}{n}$ , then

$$E(\hat{\mathbf{s}}^2) = \left(\frac{n-1}{n}\right)\sigma^2 \neq \sigma^2 = E(\mathbf{s}^2).$$

However,

$$\lim_{n \rightarrow \infty} E(\hat{\mathbf{s}}_n^2) = E(\mathbf{s}_n^2) = \sigma^2.$$

**(h)** If the population is normal, then

$$\frac{(n-1)\mathbf{s}^2}{\sigma^2} \equiv \frac{\sum_{i=1}^n (\tilde{x}_i - \bar{\mathbf{x}})^2}{\sigma^2} \sim \chi_{n-1}^2.$$

**(i)** If the population is normal, then

$$\text{Var}(\mathbf{s}^2) = \frac{2\sigma^4}{n-1}, \quad \text{Var}(\hat{\mathbf{s}}^2) = \frac{2(n-1)\sigma^4}{n^2}, \quad \text{and} \quad \text{Var}(\check{\mathbf{s}}^2) = \frac{2\sigma^4}{n}.$$

(j) If the population is normal, then

$$\frac{\bar{\mathbf{x}}_n - \boldsymbol{\mu}}{\mathbf{s} / \sqrt{n}} \sim t_{n-1}$$

(k) If  $\mathbf{s}_1^2$  and  $\mathbf{s}_2^2$  are the variances of two independent random samples of size  $n_1$  and  $n_2$  from two normal populations with variances  $\sigma_1^2$  and  $\sigma_2^2$ , respectively, then,

$$\frac{\mathbf{s}_1^2 / \sigma_1^2}{\mathbf{s}_2^2 / \sigma_2^2} = \frac{\sigma_2^2 \cdot \mathbf{s}_1^2}{\sigma_1^2 \cdot \mathbf{s}_2^2} \sim F_{n_1-1, n_2-1}.$$

In particular, if the two population variances are equal,  $\sigma_1^2 = \sigma_2^2$ , then

$$\frac{\mathbf{s}_1^2}{\mathbf{s}_2^2} \sim F_{n_1-1, n_2-1}.$$

- Assume that the population is normal. Then,

$$E(\mathbf{s}^2) = \sigma^2 \quad \text{so that } b_{\mathbf{s}^2}(\sigma^2) = 0,$$

while

$$E(\hat{\mathbf{s}}^2) = \left(\frac{n-1}{n}\right) \sigma^2 \neq \sigma^2$$

so that the bias of  $\hat{\mathbf{s}}^2$  as an estimator for  $\sigma^2$  is

$$b_{\hat{\mathbf{s}}^2}(\sigma^2) = \left(\frac{n-1}{n}\right) \sigma^2 - \sigma^2 = -\frac{\sigma^2}{n}.$$

- Moreover,

$$\text{Var}(\mathbf{s}^2) = \frac{2\sigma^4}{n-1} \quad \text{and} \quad \text{Var}(\hat{\mathbf{s}}^2) = \frac{2(n-1)\sigma^4}{n^2}.$$

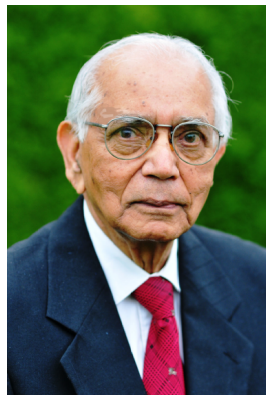
- If the population is normal,  $\mathbf{s}^2$  is an unbiased estimator for  $\sigma^2$ , while  $\hat{\mathbf{s}}^2$  is a biased estimator. However,  $\text{Var}(\hat{\mathbf{s}}^2) < \text{Var}(\mathbf{s}^2)$ . In fact, it can be proved that  $\hat{\mathbf{s}}^2$  is an estimator for  $\sigma^2$  more efficient than  $\mathbf{s}^2$ ,

$$E\left[\left(\hat{\mathbf{s}}^2 - \sigma^2\right)^2\right] < E\left[\left(\mathbf{s}^2 - \sigma^2\right)^2\right]. \quad (\text{Exercise})$$

## 8.4. The Cramér-Rao lower bound for unbiased estimators



Harald Cramér (1893 - 1985)



C. R. Rao (1920 - 2023)

- Theorem (Cramér-Rao).** Let  $\tilde{X}$  be a  $n$ -dimensional vector of random variables (not necessarily independent), the joint density of which is given by  $h(X; \theta)$ , where  $\theta$  is a  $K$ -dimensional vector of parameters in some parameter space  $\Theta$ . Let  $\hat{\theta} = g(\tilde{X})$  be an unbiased estimator for  $\theta$  with a finite covariance matrix  $\text{Var}(\hat{\theta})$ . Furthermore, assume that  $h(X; \theta)$  and  $g(X)$  satisfy

- (a)

$$\frac{\partial}{\partial \theta} \int_{\mathbb{R}^n} h dX = \int_{\mathbb{R}^n} \frac{\partial h}{\partial \theta} dX \leftarrow \text{This is a } K\text{-dimensional column vector}$$

- (b)

$$\frac{\partial}{\partial \theta} \int_{\mathbb{R}^n} g h dX = \int_{\mathbb{R}^n} g \frac{\partial h}{\partial \theta} dX \leftarrow \text{This is a } K\text{-dimensional column vector}$$

• (c)

$$\frac{\partial}{\partial \theta} \int_{\mathbb{R}^n} \frac{\partial h}{\partial \theta^T} dX = \int_{\mathbb{R}^n} \frac{\partial^2 h}{\partial \theta \partial \theta^T} dX \leftarrow \text{This is a } K \times K\text{-dimensional matrix}$$

• (d)

$$\text{Var} \left( \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right) \text{ is a positive definite } K \times K \text{ matrix.}$$

• Then, for all  $\theta \in \Theta$ ,

$$\text{Var}(\hat{\theta}) \geq \left( -\mathbb{E} \left[ \frac{\partial^2 \ln h(\tilde{X}; \theta)}{\partial \theta \partial \theta^T} \right] \right)^{-1} \equiv CR \quad (\text{CR-1})$$

or, equivalently,

$$\text{Var}(\hat{\theta}) \geq \left( \mathbb{E} \left[ \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \cdot \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta^T} \right] \right)^{-1} \equiv CR. \quad (\text{CR-2})$$

- *Note:*  $A \succeq B$  means that  $A - B$  is a positive semi-definite matrix, which implies that the elements along the diagonal of  $A - B$  are non-negative. Thus, the Cramér-Rao inequality provides a lower bound for  $\text{Var}(\hat{\theta}_k)$ , for all  $k = 1, \dots, K$ .
- **Proof.** See the handout for the case  $K = 1$ .

- **About the assumptions:**

- Assumption (d) is very mild.
- Assumptions (a), (b) and (c) mean that the operations of differentiation and integration can be interchanged. This implies, among other things, that the domain of positive density of  $\tilde{X}$  is independent of  $\theta$ . Sometimes such a domain depends on the parameters (e.g. the uniform distribution).

- The matrix

$$I(\theta) = -\mathbb{E} \left[ \frac{\partial^2 \ln h(\tilde{X}; \theta)}{\partial \theta \partial \theta^T} \right] = \mathbb{E} \left[ \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \cdot \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta^T} \right]$$

is called the Fisher information (matrix). Note that the Cramér-Rao lower bound (CR) is equal to  $[I(\theta)]^{-1}$ .

- Therefore, if  $\hat{\theta}$  is an unbiased estimator for  $\theta$  and  $\text{Var}(\hat{\theta}) = CR$ , then  $\hat{\theta}$  is the best (or the most efficient) estimator within the class of unbiased estimators for  $\theta$ . In this case, we say that the estimator  $\hat{\theta}$  is an efficient estimator for  $\theta$ .

- If the random vector  $\tilde{X}$  were discrete instead of absolutely continuous, the result would be the same. In this case  $h(X; \theta)$  would be the probability function of the random vector  $\tilde{X}$  and we should change the integrals in both the statement and the proof by sums over the range of  $\tilde{X}$ . Obviously, this range should not depend on the parameters we are estimating in order to fulfil the assumptions (a), (b) and (c).
- **Corollary.** If in the previous theorem the random vector  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  is a random sample from a population  $\tilde{x}$  having the density (or probability function)  $f(x; \theta)$ , then

$$\text{Var}(\hat{\theta}) \geq \left( -nE \left[ \frac{\partial^2 \ln f(\tilde{x}; \theta)}{\partial \theta \partial \theta^T} \right] \right)^{-1} \quad (\text{i})$$

or, equivalently,

$$\text{Var}(\hat{\theta}) \geq \left( nE \left[ \frac{\partial \ln f(\tilde{x}; \theta)}{\partial \theta} \cdot \frac{\partial \ln f(\tilde{x}; \theta)}{\partial \theta^T} \right] \right)^{-1}. \quad (\text{ii})$$

- **Proof.** See the handout for the case  $K = 1$ .

- **Single parameter case:**  $K = 1$  ( $\theta \in \mathbb{R}$ ).
- **(a)** For  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$ , not necessarily a random sample, we have

$$\text{Var}(\hat{\theta}) \geq \left( -\text{E} \left[ \frac{\partial^2 \ln h(\tilde{X}; \theta)}{\partial \theta^2} \right] \right)^{-1} \quad (\text{Theorem CR-1 for } K = 1)$$

or, equivalently,

$$\text{Var}(\hat{\theta}) \geq \left( \text{E} \left[ \left( \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right)^2 \right] \right)^{-1}. \quad (\text{Theorem CR-2 for } K = 1)$$

- **(b)** For the random sample  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$ , we have

$$\text{Var}(\hat{\theta}) \geq \left( -n\text{E} \left[ \frac{\partial^2 \ln f(\tilde{x}; \theta)}{\partial \theta^2} \right] \right)^{-1} \quad (\text{Corollary (i) for } K = 1)$$

or, equivalently,

$$\text{Var}(\hat{\theta}) \geq \left( n\text{E} \left[ \left( \frac{\partial \ln f(\tilde{x}; \theta)}{\partial \theta} \right)^2 \right] \right)^{-1}. \quad (\text{Corollary (ii) for } K = 1)$$

## 8.5. Asymptotic properties of estimators: consistent estimators

- Let  $\tilde{x}_n$  be a random variable with  $\tilde{x}_n \sim P_n$ , for  $n = 1, 2, \dots$ . The limit of the expectation (or limiting expectation) of  $\tilde{x}_n$  is

$$\lim_{n \rightarrow \infty} E(\tilde{x}_n) = \lim_{n \rightarrow \infty} \int_{\mathbb{R}} x dP_n(x),$$

whereas the asymptotic expectation is

$$AE(\tilde{x}_n) = \int_{\mathbb{R}} x d\hat{P},$$

where  $\tilde{x}_n \longrightarrow \hat{P}$  as  $n \longrightarrow \infty$ . Thus,  $\hat{P}$  is the limiting distribution of  $\{\tilde{x}_n\}_{n=1}^{\infty}$ .

- **Theorem.** If  $E(|\tilde{x}_n|^r) < M$  for all  $n$ , where  $M < \infty$  is a fixed bound, then

$$\lim_{n \rightarrow \infty} E(\tilde{x}_n^s) = AE(\tilde{x}_n^s), \quad \text{for all } s < r.$$

In particular, if  $E(\tilde{x}_n^2) < M$  for all  $n$ , where  $M < \infty$  is a fixed bound, then  $\lim_{n \rightarrow \infty} E(\tilde{x}_n) = AE(\tilde{x}_n)$ .

- Some statisticians refer to  $\lim_{n \rightarrow \infty} E(\tilde{x}_n)$  as the asymptotic expectation.
- Similarly, we can define the limit of the variance (or limiting variance) and the asymptotic variance as

$$\lim_{n \rightarrow \infty} \text{Var}(\tilde{x}_n) = \lim_{n \rightarrow \infty} E\left([\tilde{x}_n - E(\tilde{x}_n)]^2\right)$$

and

$$A\text{Var}(\tilde{x}_n) = AE\left([\tilde{x}_n - AE(\tilde{x}_n)]^2\right),$$

respectively.

- The subindex  $n$  in estimators will typically refer to the sample size.

- **Definition.** We say that the estimator  $\hat{\theta}_n$  for  $\theta \in \mathbb{R}$  is "unbiased in the limit" or "asymptotically unbiased" if

$$\lim_{n \rightarrow \infty} E(\hat{\theta}_n) = \theta \quad \text{or} \quad AE(\hat{\theta}_n) = \theta,$$

respectively.

Moreover, the "limiting bias" or the "asymptotic bias" of the estimator  $\hat{\theta}_n$  for  $\theta \in \mathbb{R}$  is

$$\lim_{n \rightarrow \infty} E(\hat{\theta}_n) - \theta \quad \text{or} \quad AE(\hat{\theta}_n) - \theta,$$

respectively.

- Definition.** An estimator  $\widehat{\theta}_n^*$  is "the best (or the most efficient) in the limit" or "asymptotically the best (or the most efficient)" in the class  $C$  of estimators for  $\theta \in \mathbb{R}$  if

$$\lim_{n \rightarrow \infty} \frac{E \left[ \left( \widehat{\theta}_n^* - \theta \right)^2 \right]}{E \left[ \left( \widehat{\theta}_n - \theta \right)^2 \right]} \leq 1 \text{ or } \frac{AE \left[ \phi(n) \left( \widehat{\theta}_n^* - \theta \right)^2 \right]}{AE \left[ \phi(n) \left( \widehat{\theta}_n - \theta \right)^2 \right]} \leq 1, \text{ for all } \widehat{\theta}_n \in C,$$

respectively.

- The function  $\phi(n)$  is selected in order to make the previous relative asymptotic efficiency well defined in  $\mathbb{R}$  (i.e., to make  $AE \left[ \phi(n) \left( \widehat{\theta}_n - \theta \right)^2 \right] \neq 0$ ). Usually, we use  $\phi(n) = n$ .

- Definition.** Let  $\hat{\theta}_n$  be an "unbiased in the limit" or an "asymptotically unbiased" estimator for  $\theta \in \mathbb{R}$ . Then,  $\hat{\theta}_n$  is an "efficient in the limit" or an "asymptotically efficient" estimator for  $\theta$  if

$$\lim_{n \rightarrow \infty} \underbrace{\left( \frac{CR_n}{\text{Var}(\hat{\theta}_n)} \right)}_{\leq 1 \text{ (Cramér-Rao)}} = 1 \quad \text{or} \quad \frac{CR^*}{\text{AVar}(\sqrt{\phi(n)}\hat{\theta}_n)} = 1,$$

respectively, where  $CR^*$  is the following "adjusted" asymptotic Cramer-Rao lower bound (or inverse of the Fisher information):

$$CR^* = \left( -\text{AE} \left[ \frac{1}{\phi(n)} \frac{\partial^2 \ln h(\tilde{X}; \theta)}{\partial \theta^2} \right] \right)^{-1} = \left( \text{AE} \left[ \left( \frac{1}{\sqrt{\phi(n)}} \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right)^2 \right] \right)^{-1}.$$

- If  $\phi(n) = n$  and  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  is a random sample from a population  $\tilde{x}$  having the density function (probability function)  $f(x; \theta)$ , then

$$CR^* = \left( -\text{AE} \left[ \frac{\partial^2 \ln f(\tilde{x}; \theta)}{\partial \theta^2} \right] \right)^{-1} = \left( \text{AE} \left[ \left( \frac{\partial \ln f(\tilde{x}; \theta)}{\partial \theta} \right)^2 \right] \right)^{-1}.$$

- *Note:* An estimator may be the best (or the most efficient) in the limit (or asymptotically) in the class of unbiased estimators in the limit (or asymptotically) even if it is not efficient in the limit (or asymptotically). In other words, the Cramér-Rao lower bound does not need to be reached in the limit (or asymptotically) by an statistic that is the best (or the most efficient) unbiased estimator in the limit (or asymptotically).

- **Definition.** An estimator  $\hat{\theta}_n$  for  $\theta \in \mathbb{R}$  is said to be (weakly) consistent if

$$\hat{\theta}_n \xrightarrow{P} \theta \quad \left( \text{or } \text{plim}_{n \rightarrow \infty} \hat{\theta}_n = \theta \right),$$

which is equivalent to  $\hat{\theta}_n \xrightarrow{d} \theta$ .

- **Definition.** An estimator  $\hat{\theta}_n$  for  $\theta \in \mathbb{R}$  is said to be strongly consistent if

$$\hat{\theta}_n \xrightarrow{a.s.} \theta.$$

- **Theorem.** If  $\lim_{n \rightarrow \infty} \text{E} \left[ \left( \hat{\theta}_n - \theta \right)^2 \right] = 0$ , then  $\hat{\theta}_n$  is a (weakly) consistent estimator for  $\theta \in \mathbb{R}$ .

- **Proof.** Obvious since  $\hat{\theta}_n \xrightarrow{m} \theta$  implies that  $\hat{\theta}_n \xrightarrow{P} \theta$ .

- Note that, since

$$\lim_{n \rightarrow \infty} \text{E} \left[ \left( \hat{\theta}_n - \theta \right)^2 \right] = \lim_{n \rightarrow \infty} \text{Var} \left( \hat{\theta}_n \right) + \lim_{n \rightarrow \infty} \left[ b_{\hat{\theta}_n}(\theta) \right]^2,$$

then  $\hat{\theta}_n \xrightarrow{m} \theta$  if and only if  $\lim_{n \rightarrow \infty} \text{Var} \left( \hat{\theta}_n \right) = \lim_{n \rightarrow \infty} b_{\hat{\theta}_n}(\theta) = 0$ .

- Let  $\bar{\mathbf{x}}_n$  and  $\mathbf{s}_n^2$  be the mean and variance of a random sample of size  $n$ .
- Since  $E(\bar{\mathbf{x}}_n) = \mu$  (or  $b_{\bar{\mathbf{x}}_n}(\mu) = 0$ ) for all  $n$ , and  $\text{Var}(\bar{\mathbf{x}}_n) = \sigma^2/n$  so that  $\lim_{n \rightarrow \infty} \text{Var}(\bar{\mathbf{x}}_n) = 0$ , the sample mean  $\bar{\mathbf{x}}_n$  is a (weakly) consistent estimator for the population mean  $\mu$ .
- In fact,  $\bar{\mathbf{x}}_n \xrightarrow{a.s.} \mu$  because of any of the strong law of large numbers.
- Assume that the population is normal. Then  $E(\mathbf{s}_n^2) = \sigma^2$  (or  $b_{\mathbf{s}_n^2}(\sigma^2) = 0$ ) for all  $n$ , and  $\text{Var}(\mathbf{s}_n^2) = 2\sigma^4/(n-1)$  so that  $\lim_{n \rightarrow \infty} \text{Var}(\mathbf{s}_n^2) = 0$ . Hence, the sample variance  $\mathbf{s}_n^2$  is a (weakly) consistent estimator for the population variance  $\sigma^2$ .

- Consider now the statistic

$$\hat{\mathbf{s}}_n^2 = \frac{\sum_{i=1}^n (\tilde{x}_i - \bar{\mathbf{x}}_n)^2}{n} = \frac{(n-1)\mathbf{s}_n^2}{n}.$$

Then,

$$E(\hat{\mathbf{s}}_n^2) = \left(\frac{n-1}{n}\right) \sigma^2 \neq \sigma^2$$

and the bias is

$$b_{\hat{\mathbf{s}}_n^2}(\sigma^2) = \left(\frac{n-1}{n}\right) \sigma^2 - \sigma^2 = -\frac{\sigma^2}{n} \longrightarrow 0, \text{ as } n \rightarrow \infty.$$

Moreover,

$$\text{Var}(\hat{\mathbf{s}}_n^2) = \frac{2(n-1)\sigma^4}{n^2} \rightarrow 0, \text{ as } n \rightarrow \infty.$$

- This means that  $\hat{\mathbf{s}}_n^2 \xrightarrow{m} \sigma^2$  and, thus,  $\hat{\mathbf{s}}_n^2$  is a (weakly) consistent estimator for  $\sigma^2$ .

## 8.6. Sufficient estimators

- The statistic  $\hat{\theta}$  is a sufficient estimator (or statistic) for the parameter  $\theta$  if it uses all the information relevant for the estimation of the population parameter  $\theta \in \Theta \subset \mathbb{R}^K$ , that is, if all the knowledge we can gain about  $\theta$  by actually specifying the individual sample values and their order can just as well be obtained by observing only the value of the statistic  $\hat{\theta}$ .
- Let  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  be a collection of random variables (a sample) whose joint distribution  $P_{\tilde{X}}(\cdot; \theta)$  depends on the parameter  $\theta$ . Then, for all values  $\hat{\theta}$  taken by  $\hat{\theta}$ , the conditional distribution of  $\tilde{X}$  given  $\hat{\theta} = \hat{\theta}$ ,  $P_{\tilde{X}|\hat{\theta}}(\cdot | \hat{\theta}; \theta)$  should not depend on the value  $\theta$ . Note that, if  $P_{\tilde{X}|\hat{\theta}}(\cdot | \hat{\theta}; \theta)$  depends on  $\theta$ , particular values of  $(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  yielding a particular value of  $\hat{\theta}$  are more likely for some values of  $\theta$  than for others, and the knowledge of these sample values will help in the estimation for  $\theta$ .

- Definition.** Let  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  be a collection of random variables (or sample) whose joint distribution  $P_{\tilde{X}}(\cdot; \theta)$  depends on the parameter  $\theta \in \Theta \subset \mathbb{R}^K$ . The statistic  $\hat{\theta} = t(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  is a sufficient estimator (or statistic) for  $\theta$  if, for all values  $\hat{\theta}$  taken by  $\hat{\theta}$ , the conditional distribution of the sample  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  given  $\hat{\theta} = \hat{\theta}$ ,  $P_{\tilde{X}|\hat{\theta}}(\cdot | \hat{\theta}; \theta) : \mathcal{B}(\mathbb{R}^n) \rightarrow [0, 1]$ , is independent of  $\theta$ .
- Note:** If  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  has a discrete (absolutely continuous) distribution, then sufficiency means that the conditional probability function (density) of  $\tilde{X}$  given  $\hat{\theta} = \hat{\theta}$ ,

$$f_{\tilde{X}|\hat{\theta}}(x_1, x_2, \dots, x_n | \hat{\theta}; \theta) = \frac{f_{(\tilde{X}, \hat{\theta})}(x_1, x_2, \dots, x_n, \hat{\theta}; \theta)}{f_{\hat{\theta}}(\hat{\theta}; \theta)},$$

for all  $\hat{\theta}$  with  $f_{\hat{\theta}}(\hat{\theta}; \theta) \neq 0$ , does not depend on the parameter  $\theta$ .

Note that

$$f_{(\tilde{X}, \hat{\theta})} \left( \underbrace{x_1, x_2, \dots, x_n}_{X \in \mathbb{R}^n}, \hat{\theta}; \theta \right) = \begin{cases} f_{\tilde{X}}(x_1, x_2, \dots, x_n; \theta) & \text{if } \hat{\theta} = t(X), \\ 0 & \text{otherwise.} \end{cases} \quad (*)$$

and, for  $\tilde{X}$  absolutely continuous,

$$f_{\hat{\theta}}(\hat{\theta}; \theta) = \int_{\mathbb{R}^n} f_{(\tilde{X}, \hat{\theta})}(x_1, x_2, \dots, x_n, \hat{\theta}; \theta) dX = \int_A f_{\tilde{X}}(x_1, x_2, \dots, x_n; \theta) dX, \\ \text{with } A = \left\{ X = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n \mid t(X) = \hat{\theta} \right\},$$

whereas, if  $\tilde{X}$  is discrete,

$$f_{\hat{\theta}}(\hat{\theta}; \theta) = \sum_{X \in \prod_{i=1}^n \tilde{x}_i(\Omega)} f_{(\tilde{X}, \hat{\theta})}(x_1, x_2, \dots, x_n, \hat{\theta}; \theta) = \sum_{X \in B} f_{\tilde{X}}(x_1, x_2, \dots, x_n; \theta), \\ \text{with } B = \left\{ X = (x_1, x_2, \dots, x_n) \in \prod_{i=1}^n \tilde{x}_i(\Omega) \mid t(X) = \hat{\theta} \right\}.$$

- Factorization theorem (Fisher-Neyman).** Let  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  be a collection of random variables (or sample) whose joint distribution is discrete (absolutely continuous) and depends on the parameter  $\theta \in \Theta \subset \mathbb{R}^K$ . Then, the statistic  $\hat{\theta} = t(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  is a sufficient estimator (or statistic) for the parameter  $\theta$  if and only if the probability function (density) of  $\tilde{X}$  can be factorized as follows:

$$f_{\tilde{X}}(x_1, x_2, \dots, x_n; \theta) = \underbrace{h(x_1, x_2, \dots, x_n)}_{\text{independent of } \theta} \cdot g(\hat{\theta}, \theta),$$

where  $\hat{\theta} = t(x_1, x_2, \dots, x_n)$ .

- **Proof.** (Sufficiency  $\implies$  Factorization). If

$$f_{\tilde{X}|\hat{\theta}}(x_1, x_2, \dots, x_n | \hat{\theta}; \theta) = \frac{f_{(\tilde{X}, \hat{\theta})}(x_1, x_2, \dots, x_n, \hat{\theta}; \theta)}{f_{\hat{\theta}}(\hat{\theta}; \theta)}$$

is independent of  $\theta$ , for all values  $\hat{\theta}$  taken by  $\hat{\theta} = t(\tilde{X})$ , then

$$\begin{aligned} & \underbrace{f_{(\tilde{X}, \hat{\theta})}(x_1, x_2, \dots, x_n, \hat{\theta}; \theta)}_{f_{\tilde{X}}(x_1, x_2, \dots, x_n; \theta) \leftarrow \text{from } (*)} \\ &= \underbrace{f_{\tilde{X}|\hat{\theta}}(x_1, x_2, \dots, x_n | \hat{\theta}; \theta)}_{\text{independent of } \theta \rightarrow m(X, \hat{\theta}) = m(X, t(X)) = h(X)} \times \underbrace{f_{\hat{\theta}}(\hat{\theta}; \theta)}_{g(\hat{\theta}, \theta)}. \end{aligned}$$

- (Factorization  $\implies$  Sufficiency). If

$$f_{\tilde{X}}(x_1, x_2, \dots, x_n; \theta) = h(x_1, x_2, \dots, x_n) \cdot g(\hat{\theta}, \theta),$$

then, for all for all values  $\hat{\theta}$  taken by  $\hat{\theta} = t(\tilde{X})$ ,

$$f_{\tilde{X}|\hat{\theta}}(x_1, x_2, \dots, x_n | \hat{\theta}; \theta) = \frac{f_{(\tilde{X}, \hat{\theta})}(x_1, x_2, \dots, x_n, \hat{\theta}; \theta) \leftarrow \text{from } (*)}{\underbrace{h(x_1, x_2, \dots, x_n) \cdot g(\hat{\theta}, \theta)}}_{f_{\hat{\theta}}(\hat{\theta}; \theta)}. \quad (**)$$

If the random vector  $\tilde{X}$  is absolutely continuous,

$$f_{\hat{\theta}}(\hat{\theta}; \theta) = \int_A \underbrace{h(x_1, x_2, \dots, x_n) \cdot g(\hat{\theta}, \theta)}_{f_{(\tilde{X}, \hat{\theta})}(x_1, x_2, \dots, x_n, \hat{\theta}; \theta)} dX = g(\hat{\theta}, \theta) \cdot \underbrace{\int_A h(x_1, x_2, \dots, x_n) dX}_{c_1(\hat{\theta})}$$

where  $A = \{X = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n \mid t(X) = \hat{\theta}\}$  or, if  $\tilde{X}$  is discrete,

$$f_{\hat{\theta}}(\hat{\theta}; \theta) = \sum_{X \in B} \underbrace{h(x_1, x_2, \dots, x_n) \cdot g(\hat{\theta}, \theta)}_{f_{(\tilde{X}, \hat{\theta})}(x_1, x_2, \dots, x_n, \hat{\theta}; \theta)} = g(\hat{\theta}, \theta) \cdot \underbrace{\sum_{X \in B} h(x_1, x_2, \dots, x_n)}_{c_2(\hat{\theta})}$$

where  $B = \{X = (x_1, x_2, \dots, x_n) \in \prod_{i=1}^n \tilde{x}_i(\Omega) \mid t(X) = \hat{\theta}\}$  and  $c_i(\hat{\theta})$ ,

$i = 1, 2$ , only depend on the value  $\hat{\theta}$  of the statistic. Therefore, from (\*\*)  
we get

$$f_{\tilde{X}|\hat{\theta}}(x_1, x_2, \dots, x_n \mid \hat{\theta}; \theta) = \frac{h(x_1, x_2, \dots, x_n)}{c_i(\hat{\theta})}, \quad i = 1, 2,$$

which is independent of  $\theta$ . Q.E.D.

## 8.7. The Rao-Blackwell theorem



C. R. Rao (1920 - 2023) - David Blackwell (1919 - 2010)

- Theorem (Rao-Blackwell).** Let  $\hat{\theta}$  be an estimator for the parameter  $\theta$ . Suppose that  $\tilde{t} = t(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  is a sufficient statistic for  $\theta$  and let  $\theta^* = E(\hat{\theta} | \tilde{t})$ . Then,

$$E[(\theta^* - \theta)^2] \leq E[(\hat{\theta} - \theta)^2]. \quad (1)$$

- Proof.** Note that, from the sufficiency of  $\tilde{t}$ ,  $\theta^* = E(\hat{\theta} | \tilde{t})$  does not depend on the parameter  $\theta$  but only on the sample  $(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  so that  $\theta^*$  is an estimator. Then,

$$\begin{aligned} E[(\theta^* - \theta)^2] &= E\left[\left[E(\hat{\theta} | \tilde{t}) - \theta\right]^2\right] = E\left[\left[E(\hat{\theta} - \theta | \tilde{t})\right]^2\right] \\ &\leq E\left[E\left[(\hat{\theta} - \theta)^2 | \tilde{t}\right]\right] = E\left[(\hat{\theta} - \theta)^2\right], \end{aligned}$$

where the inequality follows from Jensen's inequality applied to a quadratic, and thus convex, function and the last equality follows from Adam's law. *Q.E.D.*

- The previous theorem tells us that an improvement in efficiency of an estimator can be obtained by taking its conditional expectation with respect to a sufficient statistic. This process of improvement is called Rao-Blackwellization. Thus,  $\theta^* = E(\hat{\theta} | \tilde{t})$  is the Rao-Blackwellized estimator.
- In other words, the theorem says that, if an estimator is not a function of a sufficient statistic, it can be improved in terms of MSE.
- **Corollary.** Let  $\hat{\theta}$  be an unbiased estimator for the parameter  $\theta$ . Suppose that  $\tilde{t}$  is a sufficient statistic for  $\theta$  and let  $\theta^* = E(\hat{\theta} | \tilde{t})$ . Then, (a)  $\theta^*$  is an unbiased estimator for  $\theta$ , and (b)  $\text{Var}(\theta^*) \leq \text{Var}(\hat{\theta})$ .
- **Proof.** (a)  $E(\theta^*) = E[E(\hat{\theta} | \tilde{t})] = E(\hat{\theta}) = \theta$ , where the second equality follows from Adam's law and the last from the unbiasedness of  $\hat{\theta}$ . (b) From the Rao-Blackwell theorem, we know that the inequality (1) holds. Since both  $\hat{\theta}$  and  $\theta^*$  are unbiased estimators for  $\theta$ , then inequality (1) simply becomes  $\text{Var}(\theta^*) \leq \text{Var}(\hat{\theta})$ . *Q.E.D.*

## 8.8. The method of moments

- **Definition.** The  $k$ th sample moment of the value  $X = (x_1, x_2, \dots, x_n)$  taken by a random sample  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  is given by

$$m'_k(X) = \frac{\sum_{i=1}^n x_i^k}{n}.$$

- The population moments depend on the parameter  $\theta \in \Theta \subset \mathbb{R}^K$ ,  $\mu'_k(\theta)$ ,  $k = 1, 2, \dots$
- The Method of Moments (MM) consists of solving the following system of equations:

$$m'_k(X) = \mu'_k(\theta), \text{ for } k = 1, 2, \dots, K,$$

for the  $K$  parameters of the population distribution. The solution is the estimate  $\hat{\theta}_{MM} = t(X)$  and the statistic  $\hat{\theta}_{MM} = t(\tilde{X})$  is the Method of Moments estimator for  $\theta$ .

- Intuition: Find the parameters of the distribution that match the empirical moments to the population moments of the distribution.
- We could generate sample moments of order larger than  $K$ ,

$$m'_k(X) = \mu'_k(\theta), \text{ for } k = 1, 2, \dots, J, \text{ where } J \geq K.$$

However, the previous system of  $J$  equations and  $K$  unknowns does not have solution generically when  $J > K$ .

- In this case, the Generalized Method of Moments (GMM) estimate is given by

$$\hat{\theta}_{GMM} = \arg \min_{\theta \in \Theta} \{ [m'(X) - \mu'(\theta)]^T W [m'(X) - \mu'(\theta)] \},$$

where  $[m'(X) - \mu'(\theta)]^T = \left( \dots, m'_j(X) - \mu'_j(\theta), \dots \right)_{1 \times J}$  and  $W_{J \times J}$  is a weighting matrix.

- Note that if  $J = K$  and  $W = I$ , then the GMM estimator coincides with the MM estimator.

## 8.9. Maximum likelihood estimation

- Let  $X = (x_1, x_2, \dots, x_n)$  be the value of the random vector  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  having the density (probability function)  $h(X; \theta)$  where  $\theta$  is a  $K$ -dimensional vector of parameters in some parameter space  $\Theta$ . The likelihood function  $L : \Theta \rightarrow \mathbb{R}$  is given by

$$L(\theta; X) \equiv h(X; \theta).$$

- Note that  $h(X; \theta)$  is the value of the joint density (probability function) of the random vector  $\tilde{X}$  at the observed sample value.
- The maximum likelihood (ML) estimate is given by

$$\hat{\theta}_{ML} = \arg \max_{\theta \in \Theta} L(\theta; X).$$

- The solution to the previous maximization problem is the estimate  $\hat{\theta}_{ML} = t(X)$  and the statistic  $\hat{\theta}_{ML} = t(\tilde{X})$  is the Maximum Likelihood estimator for  $\theta$ .

- Intuition: Find the parameters of the distribution for which the observed data is most probable.
- If the parameter space  $\Theta$  is compact and  $L(\theta; X)$  is continuous on  $\Theta$ , then  $\hat{\theta}_{ML}$  exists.
- Sometimes, by modifying the density function on a set with zero Lebesgue measure, we can make

$$\max_{\theta \in \Theta} L(\theta; X) = \sup_{\theta \in \Theta} L(\theta; X),$$

where  $\sup_{\theta \in \Theta} L(\theta; X)$  exists if the likelihood function  $L(\cdot; X)$  is bounded and  $\Theta$  is compact.

- **Invariance property.** If  $\hat{\theta}_{ML}$  is a ML estimator for  $\theta$  and  $g : \Theta \rightarrow \Theta' \subset \mathbb{R}^K$  is a one-to-one correspondence (or bijection), then  $g(\hat{\theta}_{ML})$  is a ML estimator for  $g(\theta) \in \Theta'$ .

- Proof.** Let  $L(\theta; X)$  be the likelihood function on the parameter space  $\Theta$ ,  $\widehat{L}(\theta'; X)$  be the likelihood function on the parameter space  $\Theta'$ ,  $h(X; \theta)$  be the density (probability function) of the random vector  $\tilde{X}$ , and  $\widehat{h}(X; \theta')$  be the density function (probability function) of  $\tilde{X}$  under the reparameterization dictated by the function  $g$ , according to which  $\theta' = g(\theta)$  and  $\theta = g^{-1}(\theta')$  with  $\theta \in \Theta$  and  $\theta' \in \Theta'$ . Then, since  $h(X; \theta) = h(X; g^{-1}(\theta')) = \widehat{h}(X; \theta')$  and  $L(\theta; X) = h(X; \theta)$ , we have that

$$\widehat{L}(\theta'; X) = \widehat{h}(X; \theta') = h(X; g^{-1}(\theta')) = L(\underbrace{g^{-1}(\theta')}_{=\theta}; X).$$

Therefore,

$$\max_{\theta' \in \Theta'} \widehat{L}(\theta'; X) = \max_{g^{-1}(\theta') \in \Theta} L(g^{-1}(\theta'); X) = \max_{\theta \in \Theta} L(\theta; X).$$

The previous equality tells us that, if  $\widehat{\theta}_{ML} = \arg \max_{\theta \in \Theta} L(\theta; X)$  and

$\widehat{\theta}'_{ML} = \arg \max_{\theta' \in \Theta'} \widehat{L}(\theta'; X)$ , then  $\widehat{\theta}_{ML} = g^{-1}(\widehat{\theta}'_{ML})$  or  $\widehat{\theta}'_{ML} = g(\widehat{\theta}_{ML})$ .

*Q.E.D.*

- Example:** Let  $\hat{\theta}_{ML}$  be the Maximum Likelihood estimator for the parameter  $\theta$  of the exponential distribution. The variance of the exponential distribution is  $\sigma^2 = g(\theta) = \theta^2$  and  $g$  is a one-to-one correspondence since  $\theta > 0$ , i.e.,  $\Theta = \mathbb{R}_{++}$ . Therefore,  $\left(\hat{\theta}_{ML}\right)^2$  is the Maximum Likelihood estimator for the population variance  $\sigma^2$ .
- If  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  is a random sample from a population with density (probability function)  $f(x; \theta)$ , then

$$L(\theta; X) = h(X; \theta) = \prod_{i=1}^n f(x_i; \theta).$$

- Moreover, in this case,

$$\hat{\theta}_{ML} = \arg \max_{\theta \in \Theta} L(\theta; X) = \arg \max_{\theta \in \Theta} \ln L(\theta; X) = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \ln f(x_i; \theta).$$

## 8.10. Bayesian estimation

- $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  is a random vector (sample) and  $\tilde{w} = g(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  is a statistic, which may be a random vector.
- The random vector  $\tilde{w} = g(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  has the density (probability function)  $f_{\tilde{w}|\theta}(w|\theta)$ , where the parameter  $\theta$  is the realization of the random vector  $\theta: (\Omega, \mathcal{F}) \longrightarrow (\mathbb{R}^K, \mathcal{B}(\mathbb{R}^K))$  with  $\theta(\Omega) = \Theta \subset \mathbb{R}^K$ .
- $f_{\tilde{w}|\theta}(w|\theta)$  is also called the likelihood of  $\tilde{w}$  for  $\theta = \theta$ .
- The density (probability function) of  $\theta$  is  $f_\theta$  (this is called the prior density or prior probability function).

- The posterior density (probability function) of  $\theta$  given  $\tilde{w} = w$  is

$$f_{\theta|\tilde{w}}(\theta|w) = \frac{f_{(\theta,\tilde{w})}(\theta,w)}{f_{\tilde{w}}(w)} = \frac{f_{\theta}(\theta) f_{\tilde{w}|\theta}(w|\theta)}{f_{\tilde{w}}(w)}, \text{ for } f_{\tilde{w}}(w) \neq 0,$$

where  $f_{\theta}(\theta)$  is the prior density (probability function) of  $\theta$ ,  $f_{\tilde{w}|\theta}(w|\theta)$  is the likelihood of  $\tilde{w}$  for  $\theta = \theta$ , and

$$f_{\tilde{w}}(w) = \int_{\Theta} f_{\theta}(\theta) f_{\tilde{w}|\theta}(w|\theta) d\theta \quad (1)$$

or

$$f_{\tilde{w}}(w) = \sum_{\theta \in \Theta} f_{\theta}(\theta) f_{\tilde{w}|\theta}(w|\theta). \quad (2)$$

- Note that the formula for the posterior density/probability function,

$$f_{\theta|\tilde{w}}(\theta|w) = \frac{f_{\theta}(\theta) f_{\tilde{w}|\theta}(w|\theta)}{f_{\tilde{w}}(w)}, \text{ for } f_{\tilde{w}}(w) \neq 0,$$

also holds if  $f_{\theta}(\theta)$  is a density/probability function and the likelihood  $f_{\tilde{w}|\theta}(w|\theta)$  is a conditional probability function/density (see Exercise 13 of List 2). Then,  $f_{\tilde{w}}(w)$  would be a probability function/density obtained through the formula (1)/(2).

- We define the loss function  $L(\theta, \hat{\theta})$ . For  $K = 1$ , we usually use

$$L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2 \quad (\text{quadratic loss function})$$

or, for  $K > 1$ ,

$$L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^T W (\theta - \hat{\theta}),$$

where  $W$  is a  $K \times K$  weighting matrix.

- Then, the Bayesian estimate is given by

$$\hat{\theta}_B = \arg \min_{\hat{\theta} \in \Theta} \mathbb{E} \left[ L(\theta, \hat{\theta}) \mid \tilde{w} = w \right] = \arg \min_{\hat{\theta} \in \Theta} \int_{\Theta} L(\theta, \hat{\theta}) f_{\theta | \tilde{w}}(\theta | w) d\theta$$

or

$$\hat{\theta}_B = \arg \min_{\hat{\theta} \in \Theta} \mathbb{E} \left[ L(\theta, \hat{\theta}) \mid \tilde{w} = w \right] = \arg \min_{\hat{\theta} \in \Theta} \sum_{\theta \in \Theta} L(\theta, \hat{\theta}) f_{\theta | \tilde{w}}(\theta | w).$$

- In general,

$$\hat{\theta}_B = \arg \min_{\hat{\theta} \in \Theta} \mathbb{E} \left[ L(\theta, \hat{\theta}) \mid \tilde{w} = w \right] = \arg \min_{\hat{\theta} \in \Theta} \int_{\Theta} L(\theta, \hat{\theta}) dP_{\theta | \tilde{w}}(\theta | w).$$

- The solution is the estimate  $\hat{\theta}_B = h(w) = h(g(X)) = t(X)$  and the statistic  $\hat{\theta}_B = t(\tilde{X})$  is the Bayesian estimator for the random parameter  $\theta$  (or for the parameter value  $\theta$ ).
- Intuition: A Bayesian estimator minimizes the posterior expected value of the loss function.
- **Proposition.** Let  $K = 1$ . If  $L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$ , then  $\hat{\theta}_B = \mathbb{E}[\theta | \tilde{w} = w]$ .

- **Proof.** To solve

$$\arg \min_{\hat{\theta} \in \Theta \subset \mathbb{R}} \mathbb{E} \left[ \left( \theta - \hat{\theta} \right)^2 \mid \tilde{w} = w \right], \quad (1)$$

we compute the first derivative of  $\mathbb{E} \left[ \left( \theta - \hat{\theta} \right)^2 \mid \tilde{w} = w \right]$  w.r.t.  $\hat{\theta}$  and equate it to zero,

$$-2 \mathbb{E} \left[ \left( \theta - \hat{\theta} \right) \mid \tilde{w} = w \right] = 0.$$

Solving for  $\hat{\theta}$  in the previous expression, we get

$$\hat{\theta}_B = \mathbb{E} [\theta \mid \tilde{w} = w]. \quad (2)$$

Note that  $\mathbb{E} \left[ \left( \theta - \hat{\theta} \right)^2 \mid \tilde{w} = w \right]$  is a strictly convex function of  $\hat{\theta}$ .

To see this, we compute its second derivative w.r.t.  $\hat{\theta}$ ,

$$-2 \mathbb{E} [-1 \mid \tilde{w} = w] = (-2)(-1) = 2 > 0.$$

Therefore, the value  $\hat{\theta}_B$  obtained in (2) is a solution for (1) and, hence, it is a Bayesian estimate for  $\theta$ . *Q.E.D.*

- **Conjugate families of distributions.**
- There are cases where the prior distribution of the random vector  $\theta$  and the posterior distribution of  $\theta$  given any value of the statistic (or any value of the random sample) both belong to the same family of distributions. The prior and posterior are then called conjugate distributions for the corresponding likelihood.
- **Examples:**
- **1) Prior:**  $\lambda$  has the gamma distribution with the parameters  $\alpha$  and  $\beta$ .  
Likelihood:  $\tilde{x}$  has the Poisson distribution with the random parameter  $\lambda$ .  
Posterior of  $\lambda$  given the value of a random sample of size  $n$  from the population  $\tilde{x}$ ,  $(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n) = (x_1, x_2, \dots, x_n)$  : gamma distribution with the parameters  $\alpha + \sum_{i=1}^n x_i = \alpha + n\bar{x}$ , where  $\bar{x}$  is the value of the sample mean, and  $\frac{\beta}{1 + n\beta}$ . Thus, the family of gamma distributions is conjugate for the Poisson likelihood.

- 2) Prior:**  $\mu$  has the normal distribution with the mean  $\mu_0$  and the variance  $\sigma_0^2$  (or precision  $\tau_0 = 1/\sigma_0^2$ ). **Likelihood:**  $\tilde{x}$  has the normal distribution with the random mean  $\mu$  and the known variance  $\sigma^2$  (or precision  $\tau = 1/\sigma^2$ ), which means that the mean  $\bar{x}$  of a random sample of size  $n$  from the population  $\tilde{x}$  is normal with the random mean  $\mu$  and the known variance  $\sigma^2/n$ . **Posterior of  $\mu$  given  $\bar{x} = \bar{x}$**  (or given the sample value  $(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n) = (x_1, x_2, \dots, x_n)$ ): normal distribution with the mean  $\mu_1$  and the variance  $\sigma_1^2$  (or precision  $\tau_1 = 1/\sigma_1^2$ ), where

$$\mu_1 = \frac{n\bar{x}\sigma_0^2 + \mu_0\sigma^2}{n\sigma_0^2 + \sigma^2} = \left( \frac{\tau_0}{\tau_0 + n\tau} \right) \mu_0 + \underbrace{\left( \frac{n\tau}{\tau_0 + n\tau} \right) \bar{x}}_{= \left( \frac{\tau}{\tau_0 + n\tau} \right) \sum_{i=1}^n x_i}$$

and

$$\tau_1 = \frac{1}{\sigma_1^2} = \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} = \tau_0 + n\tau.$$

Thus, the family of normal distributions is conjugate for the normal likelihood.

- **3)** Prior:  $\theta$  has the beta distribution with the parameters  $\alpha$  and  $\beta$ .  
Likelihood:  $\tilde{x}$  has the Bernoulli distribution with the random parameter  $\theta$ . Posterior of  $\theta$  given the value of a random sample of size  $n$  from the population  $\tilde{x}$ ,  $(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n) = (x_1, x_2, \dots, x_n)$ : beta distribution with the parameters  $\sum_{i=1}^n x_i + \alpha = n\bar{x} + \alpha$  and  $n - \sum_{i=1}^n x_i + \beta = n(1 - \bar{x}) + \beta$ , where  $\bar{x}$  is the value of the sample mean (i.e., the percentage of successes). Thus, the family of beta distributions is conjugate for the Bernoulli likelihood.

OR, equivalently,

- **3')** Prior:  $\theta$  has the beta distribution with the parameters  $\alpha$  and  $\beta$ .  
Likelihood:  $\tilde{y}$  has the binomial distribution with the parameters  $\theta$  and  $n$ , where  $n$  is known. Posterior of  $\theta$  given  $\tilde{y} = y$ : beta distribution with the parameters  $y + \alpha$  and  $n - y + \beta$ . In fact, the family of beta distributions is conjugate for the binomial likelihood.

- **Sufficiency in Bayesian estimation.**

- **Definition.** Let  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  be a collection of random variables (or sample) whose joint distribution  $P_{\tilde{X}|\theta}(\cdot|\theta)$  depends on the random parameter  $\theta : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}^K, \mathcal{B}(\mathbb{R}^K))$  with  $\theta(\Omega) = \Theta \subset \mathbb{R}^K$ . The statistic  $\hat{\theta} = t(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  is a sufficient estimator (or statistic) for the random parameter  $\theta$  if, for all values  $\hat{\theta}$  taken by  $\hat{\theta}$  and all values  $\theta$  taken by  $\theta$ , the conditional distribution of the sample  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  given  $\hat{\theta} = \hat{\theta}$  and  $\theta = \theta$ ,  $P_{\tilde{X}|\hat{\theta},\theta}(\cdot|\hat{\theta},\theta) : \mathcal{B}(\mathbb{R}^n) \rightarrow [0, 1]$ , is independent of  $\theta$ .

- If  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  has a discrete (absolutely continuous) distribution, then sufficiency means that the conditional probability function (density) of  $\tilde{X}$  given  $\hat{\theta} = \hat{\theta}$  and  $\theta = \theta$ ,

$$f_{\tilde{X}|\hat{\theta},\theta}(x_1, x_2, \dots, x_n | \hat{\theta}, \theta) = \frac{f_{(\tilde{X}, \hat{\theta})}(x_1, x_2, \dots, x_n, \hat{\theta} | \theta)}{f_{\hat{\theta}|\theta}(\hat{\theta} | \theta)},$$

for all  $\hat{\theta}$  and  $\theta$  with  $f_{\hat{\theta}|\theta}(\hat{\theta} | \theta) \neq 0$ , does not depend on  $\theta$ .

- Factorization theorem.** Let  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  be a collection of random variables (or sample) whose joint distribution is discrete (absolutely continuous) and depends on the random parameter  $\theta : (\Omega, \mathcal{F}) \longrightarrow (\mathbb{R}^K, \mathcal{B}(\mathbb{R}^K))$  with  $\theta(\Omega) = \Theta \subset \mathbb{R}^K$ . Then, the statistic  $\hat{\theta} = t(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  is a sufficient estimator (or statistic) for the random parameter  $\theta$  if and only if, for all values  $\theta$  taken by  $\theta$ , the probability function (density) of  $\tilde{X}$  can be factorized as follows:

$$f_{\tilde{X}}(x_1, x_2, \dots, x_n | \theta) = h(x_1, x_2, \dots, x_n) \cdot g(\hat{\theta}, \theta),$$

where  $\hat{\theta} = t(x_1, x_2, \dots, x_n)$ .

## 8.11. Interval estimation: the pivotal method

- Let  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  be a collection of random variables (or sample) whose joint distribution depends on the parameter  $\theta \in \Theta \subset \mathbb{R}$ .
- Interval estimation is the use of a sample  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  to calculate an interval  $(a(\tilde{X}), b(\tilde{X}))$  of probable values of the unknown parameter  $\theta$ .
- **Pivotal method:** Find a real-valued random variable  $\tilde{w} = g(\tilde{X}, \theta)$ , called the pivotal variable (or pivot), such that its distribution is independent of the parameter  $\theta$  we want to estimate.
- A  $(1 - \alpha)100\%$  confidence (Borel) set  $\hat{C}_{1-\alpha}$  for  $g(\tilde{X}, \theta)$  is given by

$$P_{g(\tilde{X}, \theta)}(\hat{C}_{1-\alpha}) = P\{g(\tilde{X}, \theta) \in \hat{C}_{1-\alpha}\} = 1 - \alpha.$$

- If  $\widehat{C}_{1-\alpha}$  is an interval  $(c, d)$ , then a  $(1 - \alpha)100\%$  confidence interval for  $g(\tilde{X}, \theta)$  is given by

$$P \{c < g(\tilde{X}, \theta) < d\} = 1 - \alpha. \quad (*)$$

- Since there are many confidence intervals, we usually make

$$P \{g(\tilde{X}, \theta) \geq d\} = P \{g(\tilde{X}, \theta) \leq c\} = \frac{\alpha}{2}.$$

- Then, we solve for  $\theta$  in  $(*)$  to get a random  $(1 - \alpha)100\%$  confidence set  $H(\tilde{X}, c, d)$  for  $\theta$ ,

$$P \{\theta \in H(\tilde{X}, c, d)\} = 1 - \alpha.$$

- If the set  $H(\tilde{X}, c, d)$  is a random interval  $(a(\tilde{X}), b(\tilde{X}))$ , then

$$P \{a(\tilde{X}) < \theta < b(\tilde{X})\} = 1 - \alpha.$$

- This random interval is called the  $(1 - \alpha)100\%$  confidence interval for the parameter  $\theta$ .

- Usually, the  $(1 - \alpha)100\%$  confidence interval for  $\theta$  is presented like this:

$$a(X) < \theta < b(X),$$

where  $X$  is a particular value of the sample. Therefore, we are just providing a "realization" of the random  $(1 - \alpha)100\%$  confidence interval.

- **Examples:**

- Let  $\{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n\}$  be a random sample from a normal population with mean  $\mu$  and variance  $\sigma^2$

- **(1)** If  $\sigma^2$  is known, the random variable  $\frac{\bar{\mathbf{x}} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$ . Therefore,

$$P \left\{ -z_{\alpha/2} < \frac{\bar{\mathbf{x}} - \mu}{\sigma/\sqrt{n}} < z_{\alpha/2} \right\} = 1 - \alpha,$$

where  $z_{\alpha/2} = 2.576$  if  $1 - \alpha = 0.99$ ;  $z_{\alpha/2} = 1.96$  if  $1 - \alpha = 0.95$ ; and  $z_{\alpha/2} = 1.645$  if  $1 - \alpha = 0.9$ . Then,

$$P \left\{ \bar{\mathbf{x}} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{\mathbf{x}} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \right\} = 1 - \alpha.$$

Thus, the  $(1 - \alpha)100\%$  confidence interval for the population mean  $\mu$  is given by

$$\bar{\mathbf{x}} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{\mathbf{x}} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}},$$

where  $\bar{\mathbf{x}}$  is the value of the sample mean  $\bar{\mathbf{x}}$ .

- **(2)** If  $\sigma^2$  is unknown, the random variable  $\frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t_{n-1}$ . Therefore,

$$P \left\{ -t_{\alpha/2, n-1} < \frac{\bar{x} - \mu}{s/\sqrt{n}} < t_{\alpha/2, n-1} \right\} = 1 - \alpha.$$

Then,

$$P \left\{ \bar{x} - t_{\alpha/2, n-1} \cdot \frac{s}{\sqrt{n}} < \mu < \bar{x} + t_{\alpha/2, n-1} \cdot \frac{s}{\sqrt{n}} \right\} = 1 - \alpha.$$

Thus, the  $(1 - \alpha)100\%$  confidence interval for the population mean  $\mu$  is given by

$$\bar{x} - t_{\alpha/2, n-1} \cdot \frac{s}{\sqrt{n}} < \mu < \bar{x} + t_{\alpha/2, n-1} \cdot \frac{s}{\sqrt{n}},$$

where  $\bar{x}$  is the value of the sample mean  $\bar{x}$  and  $s$  is the value of the sample standard deviation ( $s = (s^2)^{1/2}$ ).

- **(3)** The random variable  $\frac{(n-1)s^2}{\sigma^2} \sim \chi_{n-1}^2$ . Therefore,

$$P \left\{ \chi_{1-\frac{\alpha}{2}, n-1}^2 < \frac{(n-1)s^2}{\sigma^2} < \chi_{\alpha/2, n-1}^2 \right\} = 1 - \alpha.$$

Then,

$$P \left\{ \frac{(n-1)s^2}{\chi_{\alpha/2, n-1}^2} < \sigma^2 < \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2} \right\} = 1 - \alpha.$$

Thus, the  $(1 - \alpha)100\%$  confidence interval for the population variance  $\sigma^2$  is given by

$$\frac{(n-1)s^2}{\chi_{\alpha/2, n-1}^2} < \sigma^2 < \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2},$$

where  $s^2$  is the value of the sample variance  $\mathbf{s}^2$ .

- Let  $s_1^2$  and  $s_2^2$  be the variances of independent random samples of size  $n_1$  and  $n_2$  from two normal populations with variances  $\sigma_1^2$  and  $\sigma_2^2$ , respectively.
- (4)** The random variable  $\frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2} \sim F_{n_1-1, n_2-1}$ . Therefore,

$$P \left\{ F_{1-\frac{\alpha}{2}, n_1-1, n_2-1} < \frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2} < F_{\alpha/2, n_1-1, n_2-1} \right\} = 1 - \alpha.$$

Then,

$$P \left\{ \frac{s_1^2}{s_2^2} \cdot \frac{1}{F_{\alpha/2, n_1-1, n_2-1}} < \frac{\sigma_1^2}{\sigma_2^2} < \frac{s_1^2}{s_2^2} \cdot \frac{1}{F_{1-\frac{\alpha}{2}, n_1-1, n_2-1}} \right\} = 1 - \alpha.$$

Thus, the  $(1 - \alpha)100\%$  confidence interval for  $\sigma_1^2 / \sigma_2^2$  is given by

$$\frac{s_1^2}{s_2^2} \cdot \frac{1}{F_{\alpha/2, n_1-1, n_2-1}} < \frac{\sigma_1^2}{\sigma_2^2} < \frac{s_1^2}{s_2^2} \cdot \frac{1}{F_{1-\frac{\alpha}{2}, n_1-1, n_2-1}}$$

or, since

$$F_{\alpha/2, n_2-1, n_1-1} = \frac{1}{F_{1-\frac{\alpha}{2}, n_1-1, n_2-1}}, \quad (\text{see the handout})$$

the confidence interval becomes

$$\frac{s_1^2}{s_2^2} \cdot \frac{1}{F_{\alpha/2, n_1-1, n_2-1}} < \frac{\sigma_1^2}{\sigma_2^2} < \frac{s_1^2}{s_2^2} \cdot F_{\alpha/2, n_2-1, n_1-1},$$

where  $s_1$  and  $s_2$  are the values of the standard deviations of the two random samples.

## 8.11. Confidence intervals for means, for differences between means, for proportions, for differences between proportions, for variances, and for ratios of two variances

- See the handout.

## Cramér-Rao lower bound

**Theorem (Cramér-Rao).** Let  $\tilde{X}$  be a  $n$ -dimensional vector of random variables (not necessarily independent), the joint density of which is given by  $h(X; \theta)$ , where  $\theta$  is a  $K$ -dimensional vector of parameters in some parameter space  $\Theta$ . Let  $\hat{\theta} = g(\tilde{X})$  be an unbiased estimator for  $\theta$  with a finite covariance matrix  $\text{Var}(\hat{\theta})$ . Furthermore, assume that  $h(X; \theta)$  and  $g(X)$  satisfy

(a)

$$\frac{\partial}{\partial \theta} \int_{\mathbb{R}^n} h dX = \int_{\mathbb{R}^n} \frac{\partial h}{\partial \theta} dX \leftarrow \text{This is a } K\text{-dimensional column vector}$$

(b)

$$\frac{\partial}{\partial \theta} \int_{\mathbb{R}^n} g h dX = \int_{\mathbb{R}^n} g \frac{\partial h}{\partial \theta} dX \leftarrow \text{This is a } K\text{-dimensional column vector}$$

(c)

$$\frac{\partial}{\partial \theta} \int_{\mathbb{R}^n} \frac{\partial h}{\partial \theta^T} dX = \int_{\mathbb{R}^n} \frac{\partial^2 h}{\partial \theta \partial \theta^T} dX \leftarrow \text{This is a } K \times K\text{-dimensional matrix}$$

(d)

$$\text{Var} \left( \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right) \text{ is a positive definite } K \times K \text{ matrix.}$$

Then, for all  $\theta \in \Theta$ ,

$$\text{Var}(\hat{\theta}) \geq \left( -\mathbb{E} \left[ \frac{\partial^2 \ln h(\tilde{X}; \theta)}{\partial \theta \partial \theta^T} \right] \right)^{-1} \quad (\text{CR-1})$$

or, equivalently,

$$\text{Var}(\hat{\theta}) \geq \left( \mathbb{E} \left[ \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \cdot \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta^T} \right] \right)^{-1}. \quad (\text{CR-2})$$

**Proof for the case  $K = 1$ .** We have  $1 = \int_{\mathbb{R}^n} h(X; \theta) dX$  for all  $\theta \in \Theta$ , where  $X = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$ . Therefore, differentiating both sides with respect to  $\theta$  and using assumption (a), we get

$$0 = \int_{\mathbb{R}^n} \frac{\partial h(X; \theta)}{\partial \theta} dX = \int_{\mathbb{R}^n} \frac{\partial \ln h(X; \theta)}{\partial \theta} h(X; \theta) dX \quad (1)$$

Define  $\tilde{z} = \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta}$ , then (1) means that  $E(\tilde{z}) = 0$ . Furthermore,  $\text{Var}(\tilde{z}) = E(\tilde{z}^2) = E \left[ \left( \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right)^2 \right]$ . Since  $\hat{\theta} = g(\tilde{X})$  is an unbiased estimator for  $\theta$ , we have

$$\theta = E(\hat{\theta}) \iff \theta = \int_{\mathbb{R}^n} \underbrace{g(X)}_{=\hat{\theta}} h(X; \theta) dX.$$

Therefore, differentiating both sides with respect to  $\theta$  and using assumption (b), we get

$$1 = \int_{\mathbb{R}^n} g(X) \frac{\partial h(X; \theta)}{\partial \theta} dX = \int_{\mathbb{R}^n} g(X) \frac{\partial \ln h(X; \theta)}{\partial \theta} h(X; \theta) dX. \quad (2)$$

Note that (2) means that  $E(\hat{\theta} \cdot \tilde{z}) = 1$ . Thus, we have

$$\text{Cov}(\hat{\theta} \cdot \tilde{z}) = E(\hat{\theta} \cdot \tilde{z}) - \underbrace{E(\hat{\theta})}_{=\theta} \underbrace{E(\tilde{z})}_{=0} = 1.$$

The Cauchy-Schwarz inequality implies that the absolute value of the correlation coefficient is smaller than 1 so that

$$\underbrace{[\text{Cov}(\hat{\theta} \cdot \tilde{z})]}_{=1}^2 \leq \text{Var}(\hat{\theta}) \cdot \text{Var}(\tilde{z}).$$

Note that  $[\text{Var}(\tilde{z})]^{-1}$  is well defined if assumption (d) holds. Therefore, we obtain the inequality (CR-2) since

$$\text{Var}(\hat{\theta}) \geq [\text{Var}(\tilde{z})]^{-1} = \left( E \left[ \left( \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right)^2 \right] \right)^{-1}. \quad (3)$$

To get the inequality (CR-1) observe that

$$\frac{\partial \ln h(X; \theta)}{\partial \theta} = \frac{1}{h(X; \theta)} \frac{\partial h(X; \theta)}{\partial \theta} \quad (4)$$

so that

$$\left( \frac{\partial \ln h(X; \theta)}{\partial \theta} \right)^2 = \frac{1}{[h(X; \theta)]^2} \left( \frac{\partial h(X; \theta)}{\partial \theta} \right)^2$$

and

$$E \left[ \left( \frac{\partial \ln h(X; \theta)}{\partial \theta} \right)^2 \right] = E \left[ \frac{1}{[h(X; \theta)]^2} \left( \frac{\partial h(X; \theta)}{\partial \theta} \right)^2 \right]. \quad (5)$$

Differentiating (4) with respect to  $\theta$ , we get

$$\frac{\partial^2 \ln h(X; \theta)}{\partial \theta^2} = \frac{-1}{[h(X; \theta)]^2} \left( \frac{\partial h(X; \theta)}{\partial \theta} \right)^2 + \frac{1}{h(X; \theta)} \frac{\partial^2 h(X; \theta)}{\partial \theta^2}.$$

Therefore, taking the expectation of the previous equation and using (5), we get

$$\begin{aligned}
\mathbb{E} \left[ \frac{\partial^2 \ln h(\tilde{X}; \theta)}{\partial \theta^2} \right] &= -\mathbb{E} \left[ \left( \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right)^2 \right] + \mathbb{E} \left[ \frac{1}{h(\tilde{X}; \theta)} \frac{\partial^2 h(\tilde{X}; \theta)}{\partial \theta^2} \right] \\
&= -\mathbb{E} \left[ \left( \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right)^2 \right] + \int_{\mathbb{R}^n} \frac{1}{h(X; \theta)} \frac{\partial^2 h(X; \theta)}{\partial \theta^2} h(X; \theta) dX \\
&= -\mathbb{E} \left[ \left( \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right)^2 \right] + \int_{\mathbb{R}^n} \frac{\partial^2 h(X; \theta)}{\partial \theta^2} dX, \tag{6}
\end{aligned}$$

and using assumption (c), we get

$$\int_{\mathbb{R}^n} \frac{\partial^2 h(X; \theta)}{\partial \theta^2} dX = \underbrace{\frac{\partial}{\partial \theta} \int_{\mathbb{R}^n} \frac{\partial h(X; \theta)}{\partial \theta} dX}_{=0} = 0, \tag{7}$$

where the last equality follows from differentiating both sides of the first equality in (1) with respect to  $\theta$ . Therefore, combining (6) and (7) we obtain

$$-\mathbb{E} \left[ \frac{\partial^2 \ln h(\tilde{X}; \theta)}{\partial \theta^2} \right] = \mathbb{E} \left[ \left( \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right)^2 \right]. \tag{8}$$

Plugging (8) into (3) we obtain the inequality (CR-1). *Q.E.D.*

**Corollary.** If in the previous theorem the random vector  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  is a random sample from a population  $\tilde{x}$  having the density (or probability function)  $f(x; \theta)$ , then

$$\text{Var}(\hat{\theta}) \geq \left( -n\mathbb{E} \left[ \frac{\partial^2 \ln f(\tilde{x}; \theta)}{\partial \theta \partial \theta^T} \right] \right)^{-1} \tag{i}$$

or, equivalently,

$$\text{Var}(\hat{\theta}) \geq \left( n\mathbb{E} \left[ \frac{\partial \ln f(\tilde{x}; \theta)}{\partial \theta} \cdot \frac{\partial \ln f(\tilde{x}; \theta)}{\partial \theta^T} \right] \right)^{-1}. \tag{ii}$$

**Proof for the case  $K = 1$ .** Since  $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n$  are independent, we have

$$h(X; \theta) = f(x_1; \theta) \cdot f(x_2; \theta) \cdots f(x_n; \theta),$$

so that

$$\ln h(X; \theta) = \sum_{i=1}^n \ln f(x_i; \theta),$$

which implies that

$$\frac{\partial^2 \ln h(X; \theta)}{\partial \theta^2} = \sum_{i=1}^n \frac{\partial^2 \ln f(x_i; \theta)}{\partial \theta^2},$$

and, hence,

$$\mathbb{E} \left[ \frac{\partial^2 \ln h(\tilde{X}; \theta)}{\partial \theta^2} \right] = n \mathbb{E} \left[ \frac{\partial^2 \ln f(\tilde{x}; \theta)}{\partial \theta^2} \right],$$

which combined with (CR-1) proves inequality (i) in the statement of the Corollary.

Let us prove inequality (ii). We have that

$$\begin{aligned} \mathbb{E} \left[ \left( \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right)^2 \right] &= \mathbb{E} \left[ \left( \frac{\partial [\sum_{i=1}^n \ln f(\tilde{x}_i; \theta)]}{\partial \theta} \right)^2 \right] = \mathbb{E} \left[ \left( \sum_{i=1}^n \frac{\partial \ln f(\tilde{x}_i; \theta)}{\partial \theta} \right)^2 \right] \\ &= \mathbb{E} \left[ \sum_{i=1}^n \left( \frac{\partial \ln f(\tilde{x}_i; \theta)}{\partial \theta} \right)^2 + \sum_{j \neq i} \sum_{i=1}^n \left( \frac{\partial \ln f(\tilde{x}_i; \theta)}{\partial \theta} \cdot \frac{\partial \ln f(\tilde{x}_j; \theta)}{\partial \theta} \right) \right] \\ &= \sum_{i=1}^n \left[ \mathbb{E} \left( \frac{\partial \ln f(\tilde{x}_i; \theta)}{\partial \theta} \right)^2 \right] + \sum_{j \neq i} \sum_{i=1}^n \mathbb{E} \left[ \underbrace{\frac{\partial \ln f(\tilde{x}_i; \theta)}{\partial \theta} \cdot \frac{\partial \ln f(\tilde{x}_j; \theta)}{\partial \theta}}_{\text{independent random variables}} \right] \\ &= n \mathbb{E} \left[ \left( \frac{\partial \ln f(\tilde{x}; \theta)}{\partial \theta} \right)^2 \right] + \sum_{j \neq i} \sum_{i=1}^n \underbrace{\mathbb{E} \left[ \frac{\partial \ln f(\tilde{x}_i; \theta)}{\partial \theta} \right]}_{=0} \cdot \underbrace{\mathbb{E} \left[ \frac{\partial \ln f(\tilde{x}_j; \theta)}{\partial \theta} \right]}_{=0} \end{aligned}$$

since we know from (1) that

$$0 = \mathbb{E} \left[ \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right] = \mathbb{E} \left[ \sum_{i=1}^n \frac{\partial \ln f(\tilde{x}_i; \theta)}{\partial \theta} \right] = n \mathbb{E} \left[ \frac{\partial \ln f(\tilde{x}; \theta)}{\partial \theta} \right].$$

Therefore, we have proved that

$$\mathbb{E} \left[ \left( \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right)^2 \right] = n \mathbb{E} \left[ \left( \frac{\partial \ln f(\tilde{x}; \theta)}{\partial \theta} \right)^2 \right], \quad (9)$$

which combined with (CR-2) proves inequality (ii) in the statement of the Corollary. *Q.E.D.*

**Single parameter case:**  $K = 1$  ( $\theta \in \mathbb{R}$ ).

(a) For  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$ , not necessarily a random sample, we have

$$\text{Var}(\hat{\theta}) \geq \left( -\mathbb{E} \left[ \frac{\partial^2 \ln h(\tilde{X}; \theta)}{\partial \theta^2} \right] \right)^{-1} \quad (\text{Theorem CR-1 for } K = 1)$$

or, equivalently,

$$\text{Var}(\hat{\theta}) \geq \left( \mathbb{E} \left[ \left( \frac{\partial \ln h(\tilde{X}; \theta)}{\partial \theta} \right)^2 \right] \right)^{-1}. \quad (\text{Theorem CR-2 for } K = 1)$$

(b) For the random sample  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$ , we have

$$\text{Var}(\hat{\theta}) \geq \left( -n\text{E} \left[ \frac{\partial^2 \ln f(\tilde{x}; \theta)}{\partial \theta^2} \right] \right)^{-1} \quad (\text{Corollary (i) for } K = 1)$$

or, equivalently,

$$\text{Var}(\hat{\theta}) \geq \left( n\text{E} \left[ \left( \frac{\partial \ln f(\tilde{x}; \theta)}{\partial \theta} \right)^2 \right] \right)^{-1}. \quad (\text{Corollary (ii) for } K = 1)$$

# Interval Estimation

## Confidence intervals for means

$$\frac{\bar{\mathbf{x}} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \Rightarrow P \left\{ -z_{\alpha/2} < \frac{\bar{\mathbf{x}} - \mu}{\sigma/\sqrt{n}} < z_{\alpha/2} \right\} = 1 - \alpha,$$

where  $\int_{z_{\alpha/2}}^{\infty} n(x; 0, 1)dx = \frac{\alpha}{2}$ .

Then,

$$P \left\{ \bar{\mathbf{x}} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{\mathbf{x}} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \right\} = 1 - \alpha.$$

**Proposition 1.** If  $\bar{\mathbf{x}}$  is the value of the mean of a random sample of size  $n$  from a normal population with the known variance  $\sigma^2$ , a  $(1 - \alpha)100\%$  confidence interval for the population mean  $\mu$  is given by

$$\bar{\mathbf{x}} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{\mathbf{x}} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}.$$

**Remark 1.** By virtue of the central limit theorem, this result can be also used for random samples from non-normal populations with the known variance  $\sigma^2$ , provided that  $n$  is sufficiently large ( $n \geq 30$ ). If  $\sigma$  is unknown but  $n \geq 30$ , we can replace  $\sigma$  by the value  $s$  of the sample standard deviation  $\mathbf{s} = (\mathbf{s}^2)^{1/2}$ .

-----

$$\frac{\bar{\mathbf{x}} - \mu}{\mathbf{s}/\sqrt{n}} \sim t_{n-1} \Rightarrow P \left\{ -t_{\alpha/2, n-1} < \frac{\bar{\mathbf{x}} - \mu}{\mathbf{s}/\sqrt{n}} < t_{\alpha/2, n-1} \right\} = 1 - \alpha.$$

Then,

$$P \left\{ \bar{\mathbf{x}} - t_{\alpha/2, n-1} \cdot \frac{\mathbf{s}}{\sqrt{n}} < \mu < \bar{\mathbf{x}} + t_{\alpha/2, n-1} \cdot \frac{\mathbf{s}}{\sqrt{n}} \right\} = 1 - \alpha.$$

**Proposition 2.** If  $\bar{\mathbf{x}}$  and  $s$  are the values of the mean and the standard deviation of a random sample of size  $n$  from a normal population with the unknown variance  $\sigma^2$ , a  $(1 - \alpha)100\%$  confidence interval for the population mean  $\mu$  is given by

$$\bar{\mathbf{x}} - t_{\alpha/2, n-1} \cdot \frac{\mathbf{s}}{\sqrt{n}} < \mu < \bar{\mathbf{x}} + t_{\alpha/2, n-1} \cdot \frac{\mathbf{s}}{\sqrt{n}}.$$

**Remark 2.** Note that Proposition 2 applies to small samples ( $n < 30$ ) from a normal population. For  $n \geq 30$ , the confidence interval of Proposition 2 and the one of Proposition 1 with  $\sigma$  replaced by  $s$  will generally yield nearly the same result.

### Confidence intervals for differences between means

Let  $\bar{\mathbf{x}}_1$  and  $\bar{\mathbf{x}}_2$  be the means of independent random samples of size  $n_1$  and  $n_2$  from normal populations having means  $\mu_1$  and  $\mu_2$  and variances  $\sigma_1^2$  and  $\sigma_2^2$ . Then,  $\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2 \sim N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)$ . It follows that

$$\frac{(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1).$$

From Proposition 1 we get,

**Proposition 3.** If  $\bar{\mathbf{x}}_1$  and  $\bar{\mathbf{x}}_2$  are the values of the means of independent random samples of size  $n_1$  and  $n_2$  from normal populations with the known variances  $\sigma_1^2$  and  $\sigma_2^2$ , a  $(1 - \alpha)100\%$  confidence interval for  $\mu_1 - \mu_2$  is given by

$$(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) - z_{\alpha/2} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} < \mu_1 - \mu_2 < (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) + z_{\alpha/2} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}.$$

Remark 1 also applies to Proposition 3 for large samples ( $n_1 \geq 30$  and  $n_2 \geq 30$ ).

If we had two small independent random samples from normal populations having the same variance  $\sigma^2$ , which is unknown, then we must construct the "pooled" estimator

$$\mathbf{s}_p^2 = \frac{(n_1 - 1)\mathbf{s}_1^2 + (n_2 - 1)\mathbf{s}_2^2}{n_1 + n_2 - 2}$$

which is, indeed, an unbiased estimator for  $\sigma^2$ . The independent random variables  $\frac{(n_1 - 1)\mathbf{s}_1^2}{\sigma^2}$  and  $\frac{(n_2 - 1)\mathbf{s}_2^2}{\sigma^2}$  have chi-square distributions with  $n_1 - 1$  and  $n_2 - 1$  degrees of freedom and their sum

$$\frac{(n_1 - 1)\mathbf{s}_1^2}{\sigma^2} + \frac{(n_2 - 1)\mathbf{s}_2^2}{\sigma^2} = \frac{(n_1 + n_2 - 2)\mathbf{s}_p^2}{\sigma^2} \sim \chi_{n_1 + n_2 - 2}^2.$$

Since

$$\frac{(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) - (\mu_1 - \mu_2)}{\sigma \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \quad \text{and} \quad \frac{(n_1 + n_2 - 2)\mathbf{s}_p^2}{\sigma^2}$$

are independent, it follows that

$$\frac{\frac{(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) - (\mu_1 - \mu_2)}{\sigma \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}}{\sqrt{\mathbf{s}_p^2 / \sigma^2}} = \frac{(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) - (\mu_1 - \mu_2)}{\mathbf{s}_p \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1 + n_2 - 2}.$$

From Proposition 2 we get,

**Proposition 4.** If  $\bar{x}_1$  and  $\bar{x}_2$  are the values of the means of two independent random samples of size  $n_1$  and  $n_2$  from normal populations with unknown but equal variances, a  $(1 - \alpha)100\%$  confidence interval for  $\mu_1 - \mu_2$  is given by

$$\begin{aligned} (\bar{x}_1 - \bar{x}_2) - t_{\alpha/2, n_1+n_2-2} \cdot s_p \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} &< \mu_1 - \mu_2 \\ &< (\bar{x}_1 - \bar{x}_2) + t_{\alpha/2, n_1+n_2-2} \cdot s_p \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}. \end{aligned}$$

### Confidence intervals for proportions

Let  $\tilde{x} \sim B(n, \theta)$ . Then  $\frac{\tilde{x} - n\theta}{\sqrt{n\theta(1 - \theta)}} \rightarrow N(0, 1)$  as  $n \rightarrow \infty$ .

Let  $\hat{\theta} = \frac{\tilde{x}}{n}$  be the sample proportion. Then

$$\frac{\hat{\theta} - \theta}{\sqrt{\frac{\theta(1 - \theta)}{n}}} \rightarrow N(0, 1) \text{ as } n \rightarrow \infty.$$

Using Proposition 1 and Remark 1, we get,

**Proposition 5.** An approximate  $(1 - \alpha)100\%$  confidence interval (for large  $n$ ) for the binomial parameter  $\theta$  is given by

$$\hat{\theta} - z_{\alpha/2} \cdot \sqrt{\frac{\hat{\theta}(1 - \hat{\theta})}{n}} < \theta < \hat{\theta} + z_{\alpha/2} \cdot \sqrt{\frac{\hat{\theta}(1 - \hat{\theta})}{n}},$$

where  $\hat{\theta} = \frac{x}{n}$ .

### Confidence intervals for differences between proportions

Using Proposition 3 we get,

**Proposition 6.** An approximate  $(1 - \alpha)100\%$  confidence interval (for large  $n$ ) for the difference between two binomial parameters,  $\theta_1 - \theta_2$ , is given by

$$\begin{aligned} (\hat{\theta}_1 - \hat{\theta}_2) - z_{\alpha/2} \sqrt{\frac{\hat{\theta}_1(1 - \hat{\theta}_1)}{n_1} + \frac{\hat{\theta}_2(1 - \hat{\theta}_2)}{n_2}} &< \theta_1 - \theta_2 \\ &< (\hat{\theta}_1 - \hat{\theta}_2) + z_{\alpha/2} \sqrt{\frac{\hat{\theta}_1(1 - \hat{\theta}_1)}{n_1} + \frac{\hat{\theta}_2(1 - \hat{\theta}_2)}{n_2}}. \end{aligned}$$

where  $\hat{\theta}_1 = \frac{x_1}{n_1}$  and  $\hat{\theta}_2 = \frac{x_2}{n_2}$ .

### Confidence intervals for variances

If  $\frac{(n-1)s^2}{\sigma^2} \sim \chi_{n-1}^2$ , then

$$P \left\{ \chi_{1-\frac{\alpha}{2}, n-1}^2 < \frac{(n-1)s^2}{\sigma^2} < \chi_{\alpha/2, n-1}^2 \right\} = 1 - \alpha.$$

**Proposition 7.** If  $s^2$  is the value of the variance of a random sample of size  $n$  from a normal population, a  $(1 - \alpha)100\%$  confidence interval for the population variance  $\sigma^2$  is given by

$$\frac{(n-1)s^2}{\chi_{\alpha/2, n-1}^2} < \sigma^2 < \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2}.$$

### Confidence intervals for ratios of two variances

If  $\frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2} \sim F_{n_1-1, n_2-1}$ , then

$$P \left\{ F_{1-\frac{\alpha}{2}, n_1-1, n_2-1} < \frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2} < F_{\alpha/2, n_1-1, n_2-1} \right\} = 1 - \alpha.$$

**Lemma.**  $F_{1-\alpha, \nu_1, \nu_2} = \frac{1}{F_{\alpha, \nu_2, \nu_1}}$ .

**Proof.** Exercise 8 of List 7 says that

$$\tilde{x} \sim F_{\nu_1, \nu_2} \iff \tilde{y} \equiv \frac{1}{\tilde{x}} \sim F_{\nu_2, \nu_1}.$$

Let  $P \{ \tilde{x} \geq F_{1-\alpha, \nu_1, \nu_2} \} = 1 - \alpha$ . Therefore,

$$P \left\{ \frac{1}{\tilde{y}} \geq F_{1-\alpha, \nu_1, \nu_2} \right\} = 1 - \alpha \iff P \left\{ \tilde{y} \geq \frac{1}{F_{1-\alpha, \nu_1, \nu_2}} \right\} = \alpha.$$

Since  $P \{ \tilde{y} \geq F_{\alpha, \nu_2, \nu_1} \} = \alpha$ , we arrive at the desired conclusion. *Q.E.D.*

**Proposition 8.** If  $s_1^2$  and  $s_2^2$  are the values of the variances of independent random samples of size  $n_1$  and  $n_2$  from two normal populations, a  $(1 - \alpha)100\%$  confidence interval for the ratio  $\sigma_1^2/\sigma_2^2$  of the two population variances is given by

$$\frac{s_1^2}{s_2^2} \cdot \frac{1}{F_{\alpha/2, n_1-1, n_2-1}} < \frac{\sigma_1^2}{\sigma_2^2} < \frac{s_1^2}{s_2^2} \cdot F_{\alpha/2, n_2-1, n_1-1}. \quad (\text{Check it!})$$

## Exercises. Probability and Statistics. IDEA.

### 8. Estimation

1. Show that if  $\hat{\theta}$  is an unbiased estimator for  $\theta$  and  $\text{Var}(\hat{\theta})$  does not equal 0, then  $\hat{\theta}^2$  is not an unbiased estimator for  $\theta^2$ .
2. If  $\tilde{x}_1, \tilde{x}_2, \dots,$  and  $\tilde{x}_n$  are independent Bernoulli random variables with the same parameter  $\theta$ , show that the sample proportion  $\hat{\theta} = \frac{\tilde{y}}{n}$ , where  $\tilde{y} = \sum_{i=1}^n \tilde{x}_i$ , is a minimum variance unbiased estimator for the parameter  $\theta$ . Note that  $\tilde{y}$  is a binomial random variable with the parameters  $n$  and  $\theta$ .
3. Let  $\bar{\mathbf{x}}_1$  be the mean of a random sample of size  $n$  from a normal population with the mean  $\mu$  and the variance  $\sigma_1^2$ , and  $\bar{\mathbf{x}}_2$  be the mean of a random sample of size  $n$  from a normal population with the mean  $\mu$  and the variance  $\sigma_2^2$ . If the two random samples are independent, show that
  - (a)  $w \cdot \bar{\mathbf{x}}_1 + (1 - w) \cdot \bar{\mathbf{x}}_2$ , where  $0 \leq w \leq 1$ , is an unbiased estimator for  $\mu$ ;
  - (b) the variance of the estimator is a minimum when  $w = \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2}$ .
4. If  $\bar{\mathbf{x}}_1$  and  $\bar{\mathbf{x}}_2$  are the means of independent random samples of size  $n_1$  and  $n_2$  from a normal population with the mean  $\mu$  and the variance  $\sigma^2$ , show that the variance of the unbiased estimator  $w \cdot \bar{\mathbf{x}}_1 + (1 - w) \cdot \bar{\mathbf{x}}_2$  is a minimum when  $w = \frac{n_1}{n_1 + n_2}$ .
5. Show that the mean of a random sample of size  $n$  from an exponential population  $\tilde{x}$  is a weakly consistent estimator for its parameter  $\theta$ .
6. With reference to Exercise 5 show that the sample mean is a sufficient estimator for the exponential parameter  $\theta$ . Show it in two ways: (a) without using the factorization theorem and (b) making use of the factorization theorem.
7. With reference to Exercise 2 show that the statistic  $\hat{\theta} = \frac{\tilde{y}}{n}$  is a sufficient estimator for  $\theta$ . Show it in two ways: (a) without using the factorization theorem and (b) making use of the factorization theorem.
8. Use the method of moments to find an estimator for the parameter  $\theta$  of the uniform density on the interval  $(0, \theta)$ .
9. If  $x_1, x_2, \dots,$  and  $x_n$  are the values of a random sample of size  $n$  from a population  $\tilde{x}$  having the density

$$f(x; \theta) = \begin{cases} \frac{2(\theta - x)}{\theta^2} & \text{for } 0 < x < \theta \\ 0 & \text{elsewhere,} \end{cases}$$

find an estimator for  $\theta$  by the method of moments.

10. (a) Given a random sample of size  $n$  from a normal population with the mean  $\mu$  and the variance  $\sigma^2$ , find joint maximum likelihood estimators for these two parameters.
- (b) Given a random sample of size  $n$  from a normal population with the known mean  $\mu$ , find the maximum likelihood estimator for  $\sigma$ .
11. If  $x_1, x_2, \dots$ , and  $x_n$  are the values of a random sample of size  $n$  from a geometric population  $\tilde{x}$ , find an estimate of its parameter  $\theta$  using
- (a) the method of moments;
- (b) the method of maximum likelihood.

12. Given a random sample of size  $n$  from a population  $\tilde{x}$  having the density

$$f(x; \theta) = \begin{cases} e^{-(x-\theta)} & \text{for } x > \theta \\ 0 & \text{elsewhere,} \end{cases}$$

find an estimator for the parameter  $\theta$  by the method of maximum likelihood.

13. Given independent random samples  $\{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n\}$ , and  $\{\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_n\}$  from two normal populations having the means  $\mu_1 = \alpha + \beta$  and  $\mu_2 = \alpha - \beta$  and the common variance  $\sigma^2 = 1$ , find simultaneous maximum likelihood estimators for  $\alpha$  and  $\beta$ .
14. (a) Prove that, if  $\tilde{x}$  is a binomial random variable with the parameters  $\theta$  and  $n$ , where  $n$  is known, and the prior distribution of its parameter  $\theta$  is a beta distribution with parameters  $\alpha$  and  $\beta$ , then the posterior distribution of  $\theta$  given  $\tilde{x} = x$  is a beta distribution with the parameters  $x + \alpha$  and  $n - x + \beta$ .
- (b) Under the same distributional assumptions as in (a), find the Bayesian estimate  $\hat{\theta}$  for the parameter value  $\theta$  when the loss function is  $L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$  and we observe one realization  $x$  of the random variable  $\tilde{x}$ .
- (c) Is the previous Bayesian estimator  $\hat{\theta}$  a strongly consistent estimator for the true parameter value  $\theta \in (0, 1)$ , i.e., does  $\hat{\theta} \xrightarrow{a.s.} \theta$  when  $n \rightarrow \infty$ ?

15. (a) Prove that, if  $\bar{\mathbf{x}}$  is the mean of a random sample of size  $n$  from a normal population with the known variance  $\sigma^2$ , and the prior distribution of the mean  $\boldsymbol{\mu}$  is a normal distribution with the mean  $\mu_0$  and the variance  $\sigma_0^2$ , then the posterior distribution of  $\boldsymbol{\mu}$  given  $\bar{\mathbf{x}} = \bar{x}$  is a normal distribution with the mean  $\mu_1$  and the variance  $\sigma_1^2$ , where

$$\mu_1 = \frac{n\bar{x}\sigma_0^2 + \mu_0\sigma^2}{n\sigma_0^2 + \sigma^2} \text{ and } \frac{1}{\sigma_1^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}.$$

Write the previous formulas in terms of the precisions of the previous distributions,  $\tau \equiv 1/\sigma^2$ ,  $\tau_0 \equiv 1/\sigma_0^2$ , and  $\tau_1 \equiv 1/\sigma_1^2$ .

(b) Under the same distributional assumptions as in (a), find the Bayesian estimate  $\hat{\mu}$  for the value  $\mu$  of the population mean when the loss function is  $L(\mu, \hat{\mu}) = (\mu - \hat{\mu})^2$  and we observe the value  $\bar{x}$  of the mean  $\bar{x}$  of a random sample of size  $n$ .

16. If a random sample of size  $n = 20$  from a normal population with the variance  $\sigma^2 = 225$  has the mean  $\bar{x} = 64.3$ , construct a 95% confidence interval for the population mean  $\mu$ .
17. A paint manufacturer wants to determine the average drying time of a new interior wall paint. The drying time of a given area is normally distributed. If for 12 test areas of equal size he obtained a mean drying time of 66.3 minutes and a standard deviation of 8.4 minutes, construct a 95% confidence interval for the true mean  $\mu$ .
18. Construct a 94% confidence interval for the actual difference between the average lifetimes of two kinds of light bulbs, given that a random sample of 40 light bulbs of one kind lasted on the average 418 hours of continuous use and 50 light bulbs of another kind lasted on the average 402 hours. The population standard deviations are known to be  $\sigma_1 = 26$  and  $\sigma_2 = 22$ .
19. A study has been made to compare the nicotine contents of two brands of cigarettes. Ten cigarettes of Brand *A* had an average nicotine content of 3.1 milligrams with the standard deviation of 0.5 milligram, while eight cigarettes of Brand *B* had an average nicotine content of 2.7 milligrams with a standard deviation of 0.7 milligram. Assuming that the two sets of data are random samples from normal populations with equal variances, construct a 95% confidence interval for the true difference in the average nicotine content of the two brands of cigarettes.
20. If  $x$  is a value of a random variable having an exponential distribution, find  $k$  so that the interval from 0 to  $kx$  is a  $(1 - \alpha) 100\%$  confidence interval for the parameter  $\theta$ .
21. If  $x_1$  and  $x_2$  are the values of a random sample of size 2 from a population having a uniform density on the interval  $(0, \theta)$ , find  $k$  so that

$$0 < \theta < k(x_1 + x_2)$$

is a  $(1 - \alpha) 100\%$  confidence interval for  $\theta$  with  $\alpha < 1/2$ .

22. Measurements of the blood pressure of 25 elderly women have a mean of  $\bar{x} = 140$  mm of mercury. If these data can be looked upon as a random sample from a normal population with  $\sigma = 10$  mm of mercury, construct a 95% confidence interval for the population mean  $\mu$ .
23. A study is being made to estimate the proportion of voters in a sizeable community who favor the construction of a nuclear plant. If it is found that only 140 of 400 voters selected at random favor the project, find a 95% confidence interval for the proportion of all voters in this community who favor the project.

24. If it is found that 132 of 200 voters in District  $A$  favor a given candidate for election to the United States Senate and 90 of 150 voters in District  $B$  favor the same candidate, find a 99% confidence interval for  $\theta_1 - \theta_2$ , the difference between the actual proportions of voters from the two districts who favor the candidate.
25. In 16 test runs the gasoline consumption of an experimental engine had a standard deviation of 2.2 gallons. Assume that the observed data can be looked upon as a random sample from a normal population. Construct a 99% confidence interval for  $\sigma^2$ , measuring the true variability of the gasoline consumption of this engine.
26. With reference to Exercise 19, find a 98% confidence interval for  $\frac{\sigma_1^2}{\sigma_2^2}$ .
27. Let  $\tilde{X} = \{\tilde{x}_1, \tilde{x}_2, \tilde{x}_3\}$  be a random sample from a Bernoulli population  $\tilde{x}$  with parameter  $\theta$ . Recall that  $\theta$  is the probability of success in a Bernoulli trial. Consider the statistic  $\tilde{y} = \frac{1}{6}(\tilde{x}_1 + 2\tilde{x}_2 + 3\tilde{x}_3)$ .
- (a) Find the probability function  $f(x_1, x_2, x_3, y)$  of the random vector  $(\tilde{x}_1, \tilde{x}_2, \tilde{x}_3, \tilde{y})$ .
- (b) Find the probability function  $f_{\tilde{y}}(y)$  of the random variable  $\tilde{y}$ .
- (c) Find the marginal probability function  $h(x_2, y)$  of the random vector  $(\tilde{x}_2, \tilde{y})$ .
- (d) Find the covariance between the random variables  $\tilde{x}_2$  and  $\tilde{y}$ ,  $\text{Cov}(\tilde{x}_2, \tilde{y})$ .
- (e) Find the conditional probability function of the random vector  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \tilde{x}_3)$  given  $\tilde{y} = \frac{1}{2}$  evaluated at  $X = (1, 1, 0)$ ,  $f_{\tilde{X}|\tilde{y}}(1, 1, 0 | \frac{1}{2})$ . Is the statistic  $\tilde{y}$  a sufficient estimator for  $\theta$ ?
- (f) Is the statistic  $\tilde{y}$  an unbiased estimator for  $\theta$ ? Is the statistic  $\tilde{y}$  an efficient estimator for  $\theta$ ? (i.e., is  $\tilde{y}$  an unbiased estimator whose variance equals the Cramér-Rao lower bound?) Is the statistic  $\tilde{y}$  a minimum variance unbiased estimator for  $\theta$ ? To answer this question compare the variance of  $\tilde{y}$  with the variance of the sample mean. Is the Cramér-Rao lower bound reached by the sample mean?
28. Let  $\tilde{x}_n$  be a random variable defined as follows:

$$\tilde{x}_n = \begin{cases} \tilde{z} & \text{with probability } (n-1)/n \\ n & \text{with probability } 1/n, \end{cases}$$

where  $\tilde{z} \sim N(0, 1)$ . Compute  $\text{plim}_{n \rightarrow \infty} \tilde{x}_n$ ,  $\lim_{n \rightarrow \infty} E(\tilde{x}_n)$ , and  $\text{AE}(\tilde{x}_n)$ . Note that  $\tilde{x}_n$  is a mixture.

29. Let  $\tilde{y}_n$  be a random variable defined as follows:

$$\tilde{y}_n = \begin{cases} 0 & \text{with probability } (n-1)/n \\ n^2 & \text{with probability } 1/n. \end{cases}$$

Compute  $\text{plim}_{n \rightarrow \infty} \tilde{y}_n$ ,  $\lim_{n \rightarrow \infty} E(\tilde{y}_n)$ , and  $AE(\tilde{y}_n)$ .

30. We have obtained  $x$  "successes" in  $n$  identical and independent trials.
- Find the value  $\hat{\theta}_{ML}$  of the maximum likelihood estimator of the probability  $\theta$  of success in each trial.
  - Prove that the maximum likelihood estimator  $\hat{\theta}_{ML}$  for  $\theta$  is the most efficient (or the best) estimator in the class of unbiased estimators for  $\theta$ .
  - Does the estimator  $\hat{\theta}_{ML}$  converge in mean square to  $\theta$ ,  $\hat{\theta}_{ML} \xrightarrow{m} \theta$  as  $n \rightarrow \infty$ ? Does  $\hat{\theta}_{ML} \xrightarrow{p} \theta$ ? Does  $\hat{\theta}_{ML} \xrightarrow{d} \theta$ ? Does  $\hat{\theta}_{ML} \xrightarrow{a.s.} \theta$ ?
31. In the experiment of randomly and independently choosing  $n$  real numbers from the closed interval  $[0, \beta]$ , we have obtained the numbers  $x_1, x_2, \dots, x_n$ . Let  $x^*$  be the largest of these numbers,  $x^* = \max \{x_1, x_2, \dots, x_n\}$ .
- Prove that  $x^*$  is the value of the maximum likelihood estimator for the parameter  $\beta$ .
  - Prove that the statistic  $\tilde{x}^* = \max \{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n\}$  is a weakly consistent estimator for  $\beta$ , that is,  $\tilde{x}^* \xrightarrow[n \rightarrow \infty]{p} \beta$  as  $n \rightarrow \infty$  (or  $\text{plim}_{n \rightarrow \infty} \tilde{x}^* = \beta$ ).
  - Use the method of moments to find an estimator  $\hat{\beta}_{MM}$  for the parameter  $\beta$ .
32. Let  $\{\tilde{x}_i\}_{i=1}^n$  be a random sample of size  $n$  from a population  $\tilde{x}$ . Assume that the population distribution is normal with the density

$$n(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \quad \text{for } x \in (-\infty, \infty)$$

and that the population variance is known. Prove that the maximum likelihood estimator for the population mean is also a sufficient estimator. You can use the factorization theorem for your proof.

33. Let  $\{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n\}$  be a random sample of size  $n$  from a population  $\tilde{x}$  having the Poisson distribution with a random parameter  $\lambda$ . The prior distribution of  $\lambda$  is a gamma distribution with the parameters  $\alpha$  and  $\beta$ .
- Prove that the posterior distribution of  $\lambda$  given  $(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n) = (x_1, x_2, \dots, x_n)$  is a gamma distribution with the parameters  $\alpha + \sum_{i=1}^n x_i$  and  $\frac{\beta}{1 + n\beta}$ .
  - Find the Bayesian estimate  $\hat{\lambda}_B$  for the parameter value  $\lambda$  when the loss function is quadratic,

$$L(\lambda, \hat{\lambda}) = (\lambda - \hat{\lambda})^2,$$

and we observe that  $\bar{\mathbf{x}} = \bar{x}$ , where  $\bar{\mathbf{x}}$  is the mean of a random sample with size  $n$  and  $\bar{x}$  is the realization of this sample mean.

- Compute the following biases associated with the Bayesian estimator  $\hat{\lambda}_B$ :

(i)  $E \left[ \hat{\lambda}_B - \lambda \mid \lambda = \lambda \right]$ , which is the bias when the true parameter value is  $\lambda$ , i.e., when  $\lambda = \lambda$ . Note that  $E \left[ \hat{\lambda}_B - \lambda \mid \lambda = \lambda \right] = E \left[ \hat{\lambda}_B - \lambda \mid \lambda = \lambda \right]$ .

(ii)  $E \left[ \hat{\lambda}_B - \lambda \mid \bar{x} = \bar{x} \right]$ , which is the posterior bias given a realization of the sample mean.

(iii)  $E \left[ \hat{\lambda}_B - \lambda \right]$ , which is the prior bias.

(d) Answer the following questions about the mean square error associated with the Bayesian estimator  $\hat{\lambda}_B$ :

(i) Compute  $E \left[ \left( \hat{\lambda}_B - \lambda \right)^2 \mid \lambda = \lambda \right]$ , which is the mean square error when the true parameter value is  $\lambda$ . What is the limiting bias of  $\hat{\lambda}_B$  as an estimator for the parameter value  $\lambda$ ? Is  $\hat{\lambda}_B$  a weakly consistent estimator for  $\lambda$ ? Note that  $E \left[ \left( \hat{\lambda}_B - \lambda \right)^2 \mid \lambda = \lambda \right] = E \left[ \left( \hat{\lambda}_B - \lambda \right)^2 \mid \lambda = \lambda \right]$ .

(ii) Compute  $E \left[ \left( \hat{\lambda}_B - \lambda \right)^2 \mid \bar{x} = \bar{x} \right]$ , which is the posterior mean square error given a realization of the sample mean.

(iii) Compute  $E \left[ \left( \hat{\lambda}_B - \lambda \right)^2 \right]$ , which is the prior mean square error. Does  $\text{plim}_{n \rightarrow \infty} \hat{\lambda}_B = \lambda$ ?

34. (An application of Bayesian estimation) Imagine that you have a coin that might be unbalanced. Your prior distribution of the probability  $\tilde{\theta}$  of getting a head when you flip once this coin is given by the uniform density on the interval  $(0, 1)$ . Then, you run the experiment of flipping the coin 7 times and you obtain 5 heads. What is the posterior distribution of the probability  $\tilde{\theta}$  after observing the outcome of that experiment? Find the posterior expected probability of  $\tilde{\theta}$  after observing the outcome of that experiment, i.e., compute the conditional expectation  $E \left( \tilde{\theta} \mid \tilde{x} = 5 \right)$ , where  $\tilde{x}$  is the random number of heads when flipping the coin 7 times? Compare this posterior expectation with the prior (or unconditional) expectation.

35. Let  $\{\tilde{x}_1, \tilde{x}_1, \dots, \tilde{x}_n\}$  be a random sample of size  $n$  from a normal population  $\tilde{x}$  with mean  $\mu$  and variance  $\sigma^2$ . Consider the following two statistics:

(a) the sample variance,

$$s^2 = \frac{\sum_{i=1}^n (\tilde{x}_i - \bar{x})^2}{n - 1}$$

and

(b)

$$\hat{s}^2 = \frac{\sum_{i=1}^n (\tilde{x}_i - \bar{\mathbf{x}})^2}{n},$$

where  $\bar{\mathbf{x}}$  is the sample mean.

Is the statistic given in (a) a more efficient estimator for the population variance  $\sigma^2$  than the statistic given in (b)?

36. Consider a random sample  $\{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n\}$  of size  $n$  from a normal population  $\tilde{x}$  with the known mean  $\mu$  and the unknown variance  $V > 0$ .

(a) Find the Cramér-Rao lower bound on the variance of an unbiased estimator for the population variance  $V$ .

(b) Does the variance of the sample variance  $\mathbf{s}_n^2$  reach the Cramér-Rao lower bound? Is  $\mathbf{s}_n^2$  an efficient in the limit estimator for  $V$ ?

*Note:* Recall that the parameter  $\sigma$  appearing in the normal density is the standard deviation of the population distribution and, thus,  $\sigma \equiv V^{1/2} > 0$ . Therefore, the normal density can be written as

$$\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} = \frac{1}{(2\pi)^{1/2} \cdot V^{1/2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{V}\right)^2} \equiv f(x; V), \quad -\infty < x < \infty.$$

37. Let  $\{\tilde{x}_i\}_{i=1}^n$  be a random sample of size  $n$  from a population  $\tilde{x}$ . Assume that the population distribution is Poisson with the parameter  $\lambda > 0$ .

(a) Find the the method of moments estimator for the parameter  $\lambda$ .

(b) Find the maximum likelihood estimator for the parameter  $\lambda$ .

(c) Prove that the sample mean is a sufficient estimator for the parameter  $\lambda$ . You can use the factorization theorem for your proof.

(d) Use the Cramér-Rao lower bound to check whether the sample mean is a minimum variance unbiased estimator for the parameter  $\lambda$ .

(e) The coefficient of skewness (or asymmetry) of a Poisson distribution is  $\lambda^{-1/2}$ . Find the maximum likelihood estimator for the coefficient of skewness of  $\tilde{x}$ .

38. Assume that the random variable  $\tilde{x}$  has the **Laplace distribution**, whose density is:

$$f(x; \mu, \theta) = \frac{1}{2\theta} e^{-|x-\mu|/\theta}, \quad -\infty < x < \infty.$$

Note that in the exponent of the previous density we have the negative of the absolute value  $|x - \mu|$  so that  $-|x - \mu| = -(x - \mu)$  if  $x \geq \mu$ , whereas  $-|x - \mu| = x - \mu$  if  $x \leq \mu$ .

(a) Draw the density of  $\tilde{x}$  for the particular case where  $\mu = -3$  and  $\theta = 2$ .

(b) Find the distribution function (cdf) of  $\tilde{x}$ . Is the cdf continuous? Is it differentiable everywhere? Is it twice differentiable everywhere?

(c) Find the moment generating function of  $\tilde{x}$ ? Draw it for the particular case where  $\mu = -3$  and  $\theta = 2$ .

(d) Find the mean and the variance  $\sigma^2$  of  $\tilde{x}$ .

(e) Find the density of the random variable  $\tilde{z} = |\tilde{x} - \mu|$ . Find the mean of  $\tilde{z}$ .

(f) Let  $\mu = -3$  and  $\theta = 2$ . Find the probability that  $\tilde{x} \in (-4, 1)$  given that we know that the event  $\tilde{x} \in (-5, -3)$  has occurred. *Note:* You can use for the computations of this part the distribution function you have found in part (b).

Assume for the rest of the exercise that you have a random sample  $\{\tilde{x}_i\}_{i=1}^n$  of size  $n$  from the population  $\tilde{x}$ .

(g) Find the method of moments estimator for the parameter vector  $(\mu, \theta)$ .

Assume from now on that the value of the parameter  $\mu$  is known.

(h) Find the value of the Cramér-Rao lower bound for an unbiased estimator for the parameter  $\theta$  of the distribution of  $\tilde{x}$ .

(i) Find the method of moments estimators  $\hat{\sigma}_{\text{MM}}^2$  and  $\hat{\theta}_{\text{MM}}$  for the variance  $\sigma^2$  of  $\tilde{x}$  and for the parameter  $\theta$ , respectively.

(j) Prove that the method of moments estimator  $\hat{\sigma}_{\text{MM}}^2$  is an unbiased estimator for  $\sigma^2$ . Is it consistent? (i.e., check if  $\text{plim}_{n \rightarrow \infty} \hat{\sigma}_{\text{MM}}^2 = \sigma^2$  holds).

(k) Is  $\hat{\theta}_{\text{MM}}$  an unbiased estimator for the parameter  $\theta$ ?

(l) Find the maximum likelihood estimators  $\hat{\theta}_{\text{ML}}$  and  $\hat{\sigma}_{\text{ML}}^2$  for the parameter  $\theta$  and for the variance  $\sigma^2$  of  $\tilde{x}$ , respectively.

(m) Prove that the maximum likelihood estimator  $\hat{\sigma}_{\text{ML}}^2$  is a biased estimator for the variance  $\sigma^2$  of  $\tilde{x}$ . Is it unbiased in the limit?

(n) Prove that the maximum likelihood estimator  $\hat{\theta}_{\text{ML}}$  is an unbiased estimator for  $\theta$ . Is it efficient? Is it sufficient?

39. Let  $\{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n\}$  with  $n > 1$  be a random sample from a Poisson population  $\tilde{x}$  with parameter  $\lambda$ . Consider the following unbiased estimator for  $\lambda$ :  $\hat{\lambda} = \tilde{x}_1$ . Improve this estimator in terms of mean square error using the Rao-Blackwell theorem and the result in part (c) of the previous Exercise 37.