

# Chapter 6

## Principle of Data Reduction

### 6.1 Introduction

An experimenter uses the information in a sample  $X_1, \dots, X_n$  to make inferences about an unknown parameter  $\theta$ . If the sample size  $n$  is large, then the observed sample  $x_1, \dots, x_n$  is a long list of numbers that may be hard to interpret. Any statistic,  $T(\mathbf{X})$ , defines a form of data reduction or data summary. For example, the sample mean, the sample variance, the largest observation, and the smallest observation are four statistics that might be used to summarize some key features of the sample.

The statistic summarizes the data in that, rather than reporting the entire sample  $\mathbf{x}$ , it reports only that  $T(\mathbf{x}) = t$ . For example, two samples  $\mathbf{x}$  and  $\mathbf{y}$  will be treated as equal, if  $T(\mathbf{x}) = T(\mathbf{y})$  is satisfied. The advantages and consequences of this type of data reduction are the topics of this chapter.

We study three principles of data reduction.

1. The **Sufficiency Principle** promotes a method of data reduction that does not discard information about  $\theta$  while achieving some summarization of the data.
2. The **Likelihood Principle** describes a function of the parameter determined by the observed sample, that contains all the information about  $\theta$  that is available from the sample.
3. The **Equivariance Principle** prescribes yet another method of data reduction that still preserves some important features of the model.

## 6.2 The Sufficiency Principle

**Sufficiency Principle:** If  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$ , then any inference about  $\theta$  should depend on the sample  $\mathbf{X}$  only through the value  $T(\mathbf{X})$ . That is, if  $\mathbf{x}$  and  $\mathbf{y}$  are two sample points such that  $T(\mathbf{x}) = T(\mathbf{y})$ , then the inference about  $\theta$  should be the same whether  $\mathbf{X} = \mathbf{x}$  or  $\mathbf{Y} = \mathbf{y}$  is observed.

### 6.2.1 Sufficient Statistics

**Definition 6.2.1** A statistic  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$  if the conditional distribution of the sample  $\mathbf{X}$  given the value of  $T(\mathbf{X})$  does not depend on  $\theta$ .

To use this definition to verify that a statistic  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$ , we must verify that  $P(\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x}))$  does not depend on  $\theta$ . Since  $\{\mathbf{X} = \mathbf{x}\}$  is a subset of  $\{T(\mathbf{X}) = T(\mathbf{x})\}$ ,

$$\begin{aligned} P_\theta(\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})) &= \frac{P_\theta(\mathbf{X} = \mathbf{x} \text{ and } T(\mathbf{X}) = T(\mathbf{x}))}{P_\theta(T(\mathbf{X}) = T(\mathbf{x}))} \\ &= \frac{P_\theta(\mathbf{X} = \mathbf{x})}{P_\theta(T(\mathbf{X}) = T(\mathbf{x}))} = \frac{p(\mathbf{x}|\theta)}{q(T(\mathbf{x})|\theta)}, \end{aligned}$$

where  $p(\mathbf{x}|\theta)$  is the joint pmf/pdf of the sample  $\mathbf{X}$  and  $q(t|\theta)$  is the pmf/pdf of  $T(\mathbf{X})$ . Thus,  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$  if and only if, for every  $\mathbf{x}$ , the above ratio is constant as a function of  $\theta$ .

**Theorem 6.2.1** If  $p(\mathbf{x}|\theta)$  is the joint pdf or pmf of  $\mathbf{X}$  and  $q(t|\theta)$  is the pdf or pmf of  $T(\mathbf{X})$ , then  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$  if, for every  $\mathbf{x}$  in the sample space, the ratio  $p(\mathbf{x}|\theta)/q(T(\mathbf{x})|\theta)$  is constant as a function of  $\theta$ .

**Example 6.2.1 (Binomial sufficient statistic)** Let  $X_1, \dots, X_n$  be iid Bernoulli random variables with parameter  $\theta$ ,  $0 < \theta < 1$ . Then  $T(\mathbf{X}) = X_1 + \dots + X_n$  is a sufficient statistic for  $\theta$ . Note that  $T(\mathbf{X})$  counts the number of  $X_i$ 's that equal 1, so  $T(\mathbf{X}) \sim \text{Bi}(n, \theta)$ . The ratio of pmfs is thus

$$\begin{aligned} \frac{p(\mathbf{x}|\theta)}{q(T(\mathbf{x})|\theta)} &= \frac{\prod \theta^{x_i} (1-\theta)^{1-x_i}}{\binom{n}{t} \theta^t (1-\theta)^{n-t}} \quad (\text{define } t = \sum_{i=1}^n x_i) \\ &= \frac{\theta^{\sum x_i} (1-\theta)^{\sum (1-x_i)}}{\binom{n}{t} \theta^t (1-\theta)^{n-t}} \\ &= \frac{1}{\binom{n}{t}} = \frac{1}{\binom{n}{\sum x_i}} \end{aligned}$$

Since this ratio does not depend on  $\theta$ , by Theorem 6.2.1,  $T(\mathbf{x})$  is a sufficient statistic for  $\theta$ .

**Example 6.2.2 (Normal sufficient statistic)** Let  $X_1, \dots, X_n$  be iid  $N(\mu, \sigma^2)$ , where  $\sigma^2$  is known. We wish to show that the sample mean,  $T(\mathbf{X}) = \bar{X}$ , is a sufficient statistic for  $\mu$ .

$$\begin{aligned} f(\mathbf{x}|\mu) &= \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \exp(-(x_i - \mu)^2/(2\sigma^2)) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left\{-\left[\sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2\right]/(2\sigma^2)\right\} \end{aligned}$$

Recall that the sample mean  $\bar{X} \sim N(\mu, \sigma^2/n)$ . Thus, the ratio of pdf is

$$\begin{aligned} \frac{p(\mathbf{x}|\theta)}{q(T(\mathbf{x})|\theta)} &= \frac{(2\pi\sigma^2)^{-n/2} \exp\left\{-\left[\sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2\right]/(2\sigma^2)\right\}}{(2\pi\sigma^2/n)^{-1/2} \exp\left\{-n(\bar{x} - \mu)^2/(2\sigma^2)\right\}} \\ &= n^{-1/2} (2\pi\sigma^2)^{-(n-1)/2} \exp\left\{-\sum_{i=1}^n (x_i - \bar{x})^2/(2\sigma^2)\right\}, \end{aligned}$$

which does not depend on  $\mu$ . By Theorem 6.2.1, the sample mean is a sufficient statistic for  $\mu$ .

**Example 6.2.3 (Sufficient order statistics)** Let  $X_1, \dots, X_n$  be iid from a pdf  $f$ , where we are unable to specify any more information about the pdf. It then follows that

$$f_{X_{(1)}, \dots, X_{(n)}}(\mathbf{x}) = \begin{cases} n! \prod_{i=1}^n f_X(x_i), & \text{if } x_1 < \dots < x_n \\ 0 & \text{otherwise} \end{cases}$$

where  $(x_{(1)}, \dots, x_{(n)})$  is the order statistic. By Theorem 6.2.1, the order statistic is a sufficient statistic.

Of course, this is not much of a reduction, but we should not expect more with so little information about the density  $f$ . However, even if we specify more about the density, we still may not be able to get much of a sufficient reduction. For example, suppose that  $f$  is the Cauchy pdf  $f(x|\theta) = \frac{1}{\pi(x-\theta)^2}$  or the logistic pdf  $f(x|\theta) = \frac{e^{-(x-\theta)}}{(1+e^{-(x-\theta)})^2}$ . It turns out that outside of the exponential family of distributions, it is rare to have a sufficient statistic of smaller dimension than the size of the sample, so in many cases it will turn out that the order statistics are the best that we can do.

**Theorem 6.2.2 (Factorization Theorem)** Let  $f(\mathbf{x}|\theta)$  denote the joint pdf or pmf of a sample  $\mathbf{X}$ . A statistic  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$  if and only if there exist functions  $g(t|\theta)$  and  $h(\mathbf{x})$  such that, for all sample points  $\mathbf{x}$  and all parameter points  $\theta$ ,

$$f(\mathbf{x}|\theta) = g(T(\mathbf{x})|\theta)h(\mathbf{x}). \quad (6.1)$$

PROOF: We give the proof only for discrete distributions. Suppose  $T(\mathbf{X})$  is a sufficient statistic. Choose  $g(t|\theta) = P_\theta(T(\mathbf{X}) = t)$  and  $h(\mathbf{x}) = P(\mathbf{X} = \mathbf{x}|T(\mathbf{X}) = T(\mathbf{x}))$ . Because  $T(\mathbf{X})$  is sufficient, the conditional probability  $h(\mathbf{x})$  does not depend on  $\theta$ . Thus,

$$\begin{aligned} f(\mathbf{x}|\theta) &= P_\theta(\mathbf{X} = \mathbf{x}) \\ &= P_\theta(\mathbf{X} = \mathbf{x} \text{ and } T(\mathbf{X}) = T(\mathbf{x})) \\ &= P_\theta(T(\mathbf{X}) = T(\mathbf{x}))P(\mathbf{X} = \mathbf{x}|T(\mathbf{X}) = T(\mathbf{x})) \\ &= g(T(\mathbf{x})|\theta)h(\mathbf{x}) \end{aligned}$$

So factorization (6.1) has been exhibited. Also, the last two lines above imply that  $g(T(\mathbf{x})|\theta)$  is the pmf of  $T(\mathbf{X})$ .

Now assume the factorization (6.1) exists. Let  $q(t|\theta)$  be the pmf of  $T(\mathbf{X})$ . Define  $A_{T(\mathbf{x})} = \{\mathbf{y} : T(\mathbf{y}) = T(\mathbf{x})\}$ . Then

$$\begin{aligned} \frac{f(\mathbf{x}|\theta)}{q(T(\mathbf{x})|\theta)} &= \frac{g(T(\mathbf{x})|\theta)h(\mathbf{x})}{q(T(\mathbf{x})|\theta)} \\ &= \frac{g(T(\mathbf{x})|\theta)h(\mathbf{x})}{\sum_{A_{T(\mathbf{x})}} g(T(\mathbf{y})|\theta)h(\mathbf{y})} \quad (\text{density transformation}) \\ &= \frac{g(T(\mathbf{x})|\theta)h(\mathbf{x})}{g(T(\mathbf{x})|\theta) \sum_{A_{T(\mathbf{x})}} h(\mathbf{y})} \quad (\text{since } T \text{ is a constant on } A_{T(\mathbf{x})}) \\ &= \frac{h(\mathbf{x})}{\sum_{A_{T(\mathbf{x})}} h(\mathbf{y})} \end{aligned}$$

Since the ratio does not depend on  $\theta$ , by Theorem 6.2.1,  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$ .  $\square$

To use the Factorization Theorem to find a sufficient statistic, we factor the joint pdf of the sample into two parts. One part does not depend on  $\theta$  and it constitutes the  $h(\mathbf{x})$  function. The other part depends on  $\theta$ , and usually it depends on the sample  $\mathbf{x}$  only through some function  $T(\mathbf{x})$  and this function is a sufficient statistic for  $\theta$ .

**Example 6.2.4 (Normal sufficient statistic)** Let  $X_1, \dots, X_n$  be iid  $N(\mu, \sigma^2)$ , where  $\sigma^2$  is known. The pdf can be factored as

$$f(\mathbf{x}|\mu) = (2\pi\sigma^2)^{-n/2} \exp\left\{-\sum_{i=1}^n (x_i - \bar{x})^2/(2\sigma^2)\right\} \exp\left\{-n(\bar{x} - \mu)^2/(2\sigma^2)\right\}.$$

We can define

$$g(t|\theta) = \exp\left\{-n(t - \mu)^2/(2\sigma^2)\right\}$$

by defining  $T(\mathbf{x}) = \bar{x}$ , and

$$h(\mathbf{x}) = (2\pi\sigma^2)^{-n/2} \exp\left\{-\sum_{i=1}^n (x_i - \bar{x})^2 / (2\sigma^2)\right\}.$$

Thus, by the Factorization Theorem,  $T(\mathbf{X}) = \bar{X}$  is a sufficient statistic for  $\mu$ .

**Example 6.2.5 (Uniform sufficient statistic)** Let  $X_1, \dots, X_n$  be iid observations from the discrete uniform distribution on  $1, \dots, \theta$ . That is, the unknown parameter,  $\theta$ , is a positive integer and the pmf of  $X_i$  is

$$f(x|\theta) = \begin{cases} \frac{1}{\theta} & x = 1, 2, \dots, \theta \\ 0 & \text{otherwise.} \end{cases}$$

Thus, the joint pmf of  $X_1, \dots, X_n$  is

$$f(\mathbf{x}|\theta) = \begin{cases} \theta^{-n} & x_i \in \{1, \dots, \theta\} \text{ for } i = 1, \dots, n \\ 0 & \text{otherwise.} \end{cases}$$

Let  $I_A(x)$  be the indicator function of the set  $A$ ; that is, it is equal to 1 if  $x \in A$  and equal to 0 otherwise. Let  $\mathcal{N} = \{1, 2, \dots\}$  be the set of positive integers and let  $\mathcal{N}_\theta = \{1, 2, \dots, \theta\}$ . Then the joint pmf of  $X_1, \dots, X_n$  is

$$f(\mathbf{x}|\theta) = \prod_{i=1}^n \theta^{-1} I_{\mathcal{N}_\theta}(x_i) = \theta^{-n} \prod_{i=1}^n I_{\mathcal{N}_\theta}.$$

Defining  $T(\mathbf{x}) = \max_i x_i$ , we see that

$$\prod_{i=1}^n I_{\mathcal{N}_\theta} = \left(\prod_{i=1}^n I_{\mathcal{N}}(x_i)\right) I_{\mathcal{N}_\theta}(T(\mathbf{x})).$$

Thus we have the factorization

$$f(\mathbf{x}|\theta) = \theta^{-n} I_{\mathcal{N}_\theta}(T(\mathbf{x})) \left(\prod_{i=1}^n I_{\mathcal{N}}(x_i)\right).$$

By the factorization theorem,  $T(\mathbf{X}) = \max_i X_i$  is a sufficient statistic for  $\theta$ .

**Example 6.2.6 (Normal sufficient statistic, both parameters unknown)** Assume that  $X_1, \dots, X_n$  are iid  $N(\mu, \sigma^2)$ , and that both  $\mu$  and  $\sigma^2$  are unknown so the parameter vector is  $\boldsymbol{\theta} = (\mu, \sigma^2)$ . Let  $T_1(\mathbf{x}) = \bar{x}$  and  $T_2(\mathbf{x}) = s^2 = \sum_{i=1}^n (x_i - \bar{x})^2 / (n-1)$ .

$$\begin{aligned} f(\mathbf{x}|\boldsymbol{\theta}) &= (2\pi\sigma^2)^{-n/2} \exp\left\{-\left[\sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2\right] / (2\sigma^2)\right\} \\ &= (2\pi\sigma^2)^{-n/2} \exp\left\{-(n(t_1 - \mu)^2 + (n-1)t_2) / (2\sigma^2)\right\}. \end{aligned}$$

Let  $h(\mathbf{x}) = 1$ . By the factorization theorem,  $T(\mathbf{X}) = (T_1(\mathbf{X}), T_2(\mathbf{X})) = (\bar{X}, S^2)$  is a sufficient statistic for  $(\mu, \sigma^2)$ .

The results can be generalized to the exponential family of distributions.

**Theorem 6.2.3** Let  $X_1, \dots, X_n$  be iid observations from a pdf or pmf  $f(x|\boldsymbol{\theta})$  that belongs to an exponential family given by

$$f(x|\boldsymbol{\theta}) = h(x)c(\boldsymbol{\theta}) \exp\left\{\sum_{i=1}^k w_i(\boldsymbol{\theta})t_i(x)\right\}$$

where  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d)$ ,  $d \leq k$ . Then

$$T(\mathbf{X}) = \left(\sum_{j=1}^n t_1(X_j), \dots, \sum_{j=1}^n t_k(X_j)\right)$$

is a sufficient statistic for  $\boldsymbol{\theta}$ .

### 6.2.2 Minimal Sufficient Statistics

In any problem, there are many sufficient statistics. For example,

1. It is always true that the complete sample,  $\mathbf{X}$ , is a sufficient statistic.
2. Any one-to-one function of a sufficient statistic is also a sufficient statistic.

Recall that the purpose of a sufficient statistic is to achieve data reduction without loss of information about the parameter  $\theta$ ; thus, a statistic that achieves the most data reduction while still retaining all the information about  $\theta$  might be considered preferable.

**Definition 6.2.2** A sufficient statistic  $T(\mathbf{X})$  is called a minimal sufficient statistic if, for any other sufficient statistic  $T'(\mathbf{X})$ ,  $T(\mathbf{x})$  is a function of  $T'(\mathbf{x})$ .

To say that  $T(\mathbf{x})$  is a function of  $T'(\mathbf{x})$  simply means that if  $T'(\mathbf{x}) = T'(\mathbf{y})$ , then  $T(\mathbf{x}) = T(\mathbf{y})$ . Let  $\mathcal{T} = \{t : t = T(\mathbf{x}) \text{ for some } \mathbf{x} \in \mathcal{X}\}$  be the image of  $\mathcal{X}$  under  $T(\mathbf{x})$ . Then  $T\mathcal{X}$  partitions the sample space into sets  $A_t, t \in \mathcal{T}$ , defined by  $A_t = \{\mathbf{x} : T(\mathbf{x}) = t\}$ . If  $\{B_{t'} : t' \in \mathcal{T}'\}$  are the partition sets for  $T'(\mathbf{x})$  and  $\{A_t : t \in \mathcal{T}\}$  are the partition sets for  $T(\mathbf{x})$ , then Definition 6.2.2 states that every  $B_{t'}$  is a subset of some  $A_t$ . Thus, the partition associates with a minimal sufficient statistic, is the coarsest possible partition for a sufficient statistic, and a minimal sufficient statistic achieves the greatest possible data reduction for a sufficient statistic.

**Example 6.2.7 (Two Normal sufficient statistics)** Suppose that  $X_1, \dots, X_n$  are observations from  $N(\mu, \sigma^2)$  with  $\sigma^2$  known. We know that  $T(\mathbf{X}) = \bar{X}$  is a sufficient statistic for  $\mu$ . The factorization theorem shows that  $T'(\mathbf{X}) = (\bar{X}, S^2)$  is also a sufficient statistic for  $\mu$ .  $T(\mathbf{X})$  can be written as a function of  $T'(\mathbf{X})$  by defining the function  $r(a, b) = a$ . Then  $T(\mathbf{x}) = \bar{x} = r(\bar{x}, S^2) = r(T'(\mathbf{x}))$ .

**Theorem 6.2.4** let  $f(\mathbf{x}|\theta)$  be the pmf or pdf of a sample  $\mathbf{X}$ . Suppose there exists a function  $T(\mathbf{x})$  such that, for every two sample points  $\mathbf{x}$  and  $\mathbf{y}$ , the ratio  $f(\mathbf{x}|\theta)/f(\mathbf{y}|\theta)$  is a constant as a function of  $\theta$  if and only if  $T(\mathbf{x}) = T(\mathbf{y})$ . Then  $T(\mathbf{x})$  is a minimal sufficient statistic for  $\theta$ .

PROOF: To simplify the proof, we assume  $f(\mathbf{x}|\theta) > 0$  for all  $\mathbf{x} \in \mathcal{X}$  and  $\theta$ . First we show that  $T(\mathbf{X})$  is a sufficient statistic. Let  $\mathcal{T} = \{t : t = T(\mathbf{x}) \text{ for some } \mathbf{x} \in \mathcal{X}\}$  be the image of  $\mathcal{X}$  under  $T(\mathbf{x})$ . Define the partition sets induced by  $T(\mathbf{x})$  as  $A_t = \{\mathbf{x} : T(\mathbf{x}) = t\}$ . For each  $A_t$ , choose and fix one element  $\mathbf{x}_t \in A_t$ . For any  $\mathbf{x} \in \mathcal{T}$ ,  $\mathbf{x}_{T(\mathbf{x})}$  is the fixed element that is in the same set,  $A_t$ , as  $\mathbf{x}$ . Since  $\mathbf{x}$  and  $\mathbf{x}_{T(\mathbf{x})}$  are in the same set  $A_t$ ,  $T(\mathbf{x}) = T(\mathbf{x}_{T(\mathbf{x})})$  and, hence,  $f(\mathbf{x}|\theta)/f(\mathbf{x}_{T(\mathbf{x})}|\theta)$  is constant as a function of  $\theta$ . Thus, we can define a function on  $\mathcal{X}$  by  $h(\mathbf{x}) = f(\mathbf{x}|\theta)/f(\mathbf{x}_{T(\mathbf{x})}|\theta)$  and  $h$  does not depend on  $\theta$ . Define a function on  $\mathcal{T}$  by  $g(t|\theta) = f(\mathbf{x}_t|\theta)$ . Then it can be seen that

$$f(\mathbf{x}|\theta) = \frac{f(\mathbf{x}_{T(\mathbf{x})}|\theta)f(\mathbf{x}|\theta)}{f(\mathbf{x}_{T(\mathbf{x})}|\theta)} = g(T(\mathbf{x})|\theta)h(\mathbf{x})$$

and, by the factorization theorem,  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$ .

Now to show that  $T(\mathbf{X})$  is minimal, let  $T'(\mathbf{X})$  be any other sufficient statistic. By the factorization theorem, there exist function  $g'$  and  $h'$  such that  $f(\mathbf{x}|\theta) = g'(T'(\mathbf{x})|\theta)h'(\mathbf{x})$ . Let  $\mathbf{x}$  and  $\mathbf{y}$  be any two sample points with  $T'(\mathbf{x}) = T'(\mathbf{y})$ . Then

$$\frac{f(\mathbf{x}|\theta)}{f(\mathbf{y}|\theta)} = \frac{g'(T'(\mathbf{x})|\theta)h'(\mathbf{x})}{g'(T'(\mathbf{y})|\theta)h'(\mathbf{y})} = \frac{h'(\mathbf{x})}{h'(\mathbf{y})}.$$

Since this ratio does not depend on  $\theta$ , the assumptions of the theorem imply that  $T(\mathbf{x}) = T(\mathbf{y})$ . Thus,  $T(\mathbf{x})$  is a function of  $T'(\mathbf{x})$  and  $T(\mathbf{x})$  is minimal.  $\square$

**Example 6.2.8 (Normal minimal sufficient statistic)** Let  $X_1, \dots, X_n$  be iid  $N(\mu, \sigma^2)$ , both  $\mu$  and  $\sigma^2$  unknown. Let  $\mathbf{x}$  and  $\mathbf{y}$  denote two sample points, and let  $(\bar{x}, S_{\mathbf{x}}^2)$  and  $(\bar{y}, S_{\mathbf{y}}^2)$  be the sample means and variances corresponding to the  $\mathbf{x}$  and  $\mathbf{y}$  samples, respectively. Then the ratio of densities is

$$\frac{f(\mathbf{x}|\mu, \sigma^2)}{f(\mathbf{y}|\mu, \sigma^2)} = \exp\{[-n(\bar{x}^2 - \bar{y}^2) + 2n\mu(\bar{x} - \bar{y}) - (n-1)(s_{\mathbf{x}}^2 - s_{\mathbf{y}}^2)]/(2\sigma^2)\}.$$

This ratio will be constant as a function of  $\mu$  and  $\sigma^2$  if and only if  $\bar{x} = \bar{y}$  and  $s_{\mathbf{x}}^2 = s_{\mathbf{y}}^2$ . Thus, by Theorem 6.2.4,  $(\bar{X}, S^2)$  is a minimal sufficient statistic for  $(\mu, \sigma^2)$ .

**Example 6.2.9 (Uniform minimal sufficient statistic)** Suppose  $X_1, \dots, X_n$  are iid uniform observations on the interval  $(\theta, \theta + 1)$ ,  $-\infty < \theta < \infty$ . Then the joint pdf of  $X$  is

$$f(\mathbf{x}|\theta) = \begin{cases} 1 & \theta < x_i < \theta + 1, i = 1, \dots, n, \\ 0 & \text{otherwise,} \end{cases}$$

which can be written as

$$f(\mathbf{x}|\theta) = \begin{cases} 1 & \max_i x_i - 1 < \theta < \min_i x_i \\ 0 & \text{otherwise} \end{cases}$$

Thus, for two sample points  $\mathbf{x}$  and  $\mathbf{y}$ , the numerator and denominator of the ratio  $f(\mathbf{x}|\theta)/f(\mathbf{y}|\theta)$  will be positive for the same values of  $\theta$  if and only if  $\min_i x_i = \min_i y_i$  and  $\max_i x_i = \max_i y_i$ . Thus, we have that  $T(\mathbf{X}) = (X_{(1)}, X_{(n)})$  is a minimal sufficient statistic.

This example shows the dimension of a minimal sufficient statistic may not match the dimension of the parameter. A minimal sufficient statistic is not unique. Any one-to-one function of a minimal sufficient statistic is also a minimal sufficient statistic. For example,  $T'(\mathbf{X}) = (X_{(n)} - X_{(1)}, (X_{(n)} + X_{(1)})/2)$  is also a minimal sufficient statistic in Example 6.2.9, and  $T'(\mathbf{X}) = (\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2)$  is also a minimal sufficient statistic in Example 6.2.8.

### 6.2.3 Ancillary Statistics

**Definition 6.2.3** A statistic  $S(\mathbf{X})$  whose distribution does not depend on the parameter  $\theta$  is called an ancillary statistic.

Alone, an ancillary statistic contains no information about  $\theta$ . An ancillary statistic is an observation on a random variable whose distribution is fixed and known, unrelated to  $\theta$ . Paradoxically, an ancillary statistic, when used in conjunction with other statistics, sometimes does contain valuable information for inferences about  $\theta$ .

**Example 6.2.10 (Uniform ancillary statistic)** Let  $X_1, \dots, X_n$  be iid uniform observations on the interval  $(\theta, \theta + 1)$ ,  $-\infty < \theta < \infty$ . The range statistic  $R = X_{(n)} - X_{(1)}$  is an ancillary statistic.

Recall that the cdf of each  $X_i$  is

$$F(x|\theta) = \begin{cases} 0 & x \leq \theta \\ x - \theta & \theta < x < \theta + 1 \\ 1 & \theta + 1 \leq x. \end{cases}$$

Thus, the joint pdf of  $X_{(1)}$  and  $X_{(n)}$  is

$$g(x_{(1)}, x_{(n)}|\theta) = \begin{cases} n(n-1)(x_{(n)} - x_{(1)})^{n-2} & \theta < x_{(1)} < x_{(n)} < \theta + 1 \\ 0 & \text{otherwise} \end{cases}$$

Making the transformation  $R = X_{(n)} - X_{(1)}$  and  $M = (X_{(n)} + X_{(1)})/2$ , we see the joint pdf of  $R$  and  $M$  is

$$h(r, m|\theta) = \begin{cases} n(n-1)r^{n-2} & 0 < r < 1, \theta + r/2 < m < \theta + 1 - r/2 \\ 0 & \text{otherwise} \end{cases}$$

Thus, the pdf for  $R$  is

$$h(r|\theta) = \int_{\theta+r/2}^{\theta+1-r/2} n(n-1)r^{n-2} dm = n(n-1)r^{n-2}(1-r), \quad 0 < r < 1.$$

This is a beta pdf with  $\alpha = n - 1$  and  $\beta = 2$ . Thus, the distribution of  $R$  does not depend on  $\theta$ , and  $R$  is ancillary.

**Example 6.2.11 (Location and scale family ancillary statistic)** Let  $X_1, \dots, X_n$  be iid observations from a location family with cdf  $F(x - \theta)$ ,  $-\infty < \theta < \infty$ . Show that the range statistic  $R = X_{(n)} - X_{(1)}$  is an ancillary statistic.

Let  $X_1, \dots, X_n$  be iid observations from a scale family with cdf  $F(x/\sigma)$ ,  $\sigma > 0$ . Show that any statistic that depends on the sample only through the  $n - 1$  values  $X_1/X_n, \dots, X_{n-1}/X_n$  is an ancillary statistic.

#### 6.2.4 Sufficient, Ancillary, and Complete Statistics

A minimal sufficient statistic is a statistic that has achieved the maximal amount of data reduction possible while still retaining all the information about the parameter  $\theta$ . Intuitively, a minimal sufficient statistic eliminates all the extraneous information in the sample, retaining only that piece with information about  $\theta$ . Since the distribution of an ancillary statistic does not depend on  $\theta$ , it

might be suspected that a minimal sufficient statistic is unrelated to an ancillary statistic. However, this is not necessarily the case. Recall Example 6.2.9 in which  $X_1, \dots, X_n$  were iid obs from a  $\text{uniform}(\theta, \theta + 1)$  distribution. We have pointed out that the statistic  $(X_{(n)} - X_{(1)}, (X_{(n)} + X_{(1)})/2)$  is a minimal sufficient statistic, and in Example 6.2.10 we showed that  $X_{(n)} - X_{(1)}$  is an ancillary statistic. Thus, in this case, the ancillary statistic is an important component of the minimal sufficient statistic. Certainly, the ancillary statistic and the minimal sufficient statistic are not independent.

The following example shows that an ancillary statistic can sometimes give important information for inference about  $\theta$ .

**Example 6.2.12 (Ancillary precision)** *Let  $X_1$  and  $X_2$  be iid obs from the discrete distribution that satisfies*

$$P(X = \theta) = P(X = \theta + 1) = P(\theta + 2) = 1/3,$$

where  $\theta$ , the unknown parameter, is any integer. It can be shown with an argument similar to that in Example 6.2.9 that  $(R, M)$ , where  $R = X_{(n)} - X_{(1)}$  and  $M = (X_{(n)} + X_{(1)})/2$ , is a minimal sufficient statistic. To see how  $R$  might give information about  $\theta$ , consider a sample point  $(r, m)$ , where  $m$  is an integer. First we consider only  $m$ ,  $\theta$  must be one of three values, either  $\theta = m$  or  $\theta = m - 1$  or  $\theta = m - 2$ . With only the information that  $M = m$ , all three  $\theta$  values are possible values. But now suppose we get the information that  $R = 2$ . Then it must be the case that  $X_{(1)} = m - 1$  and  $X_{(2)} = m + 1$ , and the only possible value for  $\theta$  is  $\theta = m - 1$ . Thus, the knowledge of the value of the ancillary statistic  $R$  has increased our knowledge about  $\theta$ . Of course, the knowledge of  $R$  alone would give us no information about  $\theta$ .

For many important situations, however, our intuition that a minimal sufficient statistic is independent of any ancillary statistic is correct. A description of situations in which this occurs relies on the next definition.

**Definition 6.2.4** *Let  $f(t|\theta)$  be a family of pdfs or pmfs for a statistic  $T(\mathbf{X})$ . The family of probability distributions is called complete if  $E_\theta g(T) = 0$  for all  $\theta$  implies  $P_\theta(g(T) = 0) = 1$  for all  $\theta$ . Equivalently,  $T(\mathbf{X})$  is called a complete statistic.*

**Example 6.2.13** (Binomial complete sufficient statistic) *Suppose that  $T$  has a binomial  $(n, p)$  distribution,  $0 < p < 1$ . Let  $g$  be a function such that  $E_p g(T) = 0$ . Then*

$$0 = E_p g(T) = \sum_{t=0}^n g(t) \binom{n}{t} p^t (1-p)^{n-t} = (1-p)^n \sum_{t=0}^n g(t) \binom{n}{t} \left(\frac{p}{1-p}\right)^t$$

for all  $p$ ,  $0 < p < 1$ . Let  $r = p/(1-p)$ ,  $0 < r < \infty$ . The last expression is a polynomial of degree  $n$  in  $r$ , where the coefficient of  $r^t$  is  $g(t) \binom{n}{t}$ . For the polynomial to be 0 for all  $r$ , each coefficient must be 0. Thus,  $g(t) = 0$  for  $t = 0, 1, \dots, n$ . Since  $T$  takes on the values  $0, 1, \dots, n$  with probability 1, this yields that  $P_p(g(T) = 0) = 1$  for all  $p$ . Hence,  $T$  is a complete statistic.

**Example 6.2.14** *Let  $X_1, \dots, X_n$  be iid uniform  $(0, \theta)$  observations,  $0 < \theta < \infty$ . Show that  $T(\mathbf{X}) = \max_i X_i$  is a complete statistic.*

Using an argument similar to that used in one of previous examples, we can see that  $T(\mathbf{X}) = \max_i X_i$  is a sufficient statistic and the pdf of  $T(\mathbf{X})$  is

$$f(t|\theta) = \begin{cases} nt^{n-1}\theta^{-n} & 0 < t < \theta \\ 0 & \text{otherwise.} \end{cases}$$

Suppose  $g(t)$  is a function satisfying  $E_\theta g(T) = 0$  for all  $\theta$ . Since  $E_\theta g(T)$  is a constant of  $\theta$ , its derivative with respect to  $\theta$  is 0. Thus we have that

$$\begin{aligned} 0 &= \frac{d}{d\theta} E_\theta g(T) = \frac{d}{d\theta} \int_0^\theta g(t) n t^{n-1} \theta^{-n} dt \\ &= (\theta^{-n}) \frac{d}{d\theta} \int_0^\theta n g(t) t^{n-1} dt + \left(\frac{d}{d\theta} \theta^{-n}\right) \int_0^\theta n g(t) t^{n-1} dt \\ &= \theta^{-n} n g(\theta) \theta^{n-1} + 0 \\ &= \theta^{-1} n g(\theta) \end{aligned}$$

Since  $\theta^{-1} n \neq 0$ , it must be  $g(\theta) = 0$ . This is true for every  $\theta > 0$ , hence,  $T$  is a complete statistic.

**Theorem 6.2.5 (Basu's Theorem)** *If  $T(\mathbf{X})$  is a complete and minimal sufficient statistic, then  $T(\mathbf{X})$  is independent of every ancillary statistic.*

PROOF: We give the proof only for discrete distributions. Let  $S(X)$  be any ancillary statistic. Then  $P(S(X) = s)$  does not depend on  $\theta$  since  $S(X)$  is ancillary. Also the conditional probability,

$$P(S(X) = s | T(\mathbf{X}) = t) = P(\mathbf{X} \in \{\mathbf{x} : S(\mathbf{x}) = s\} | T(\mathbf{X}) = t),$$

does not depend on  $\theta$  because  $T(\mathbf{X})$  is a sufficient statistic (recall the definition!). Thus, to show that  $S(\mathbf{X})$  and  $T(\mathbf{X})$  are independent, it suffices to show that

$$P(S(\mathbf{X}) = s | T(\mathbf{X}) = t) = P(S(\mathbf{X}) = s) \quad (6.2)$$

for all possible values  $t \in \mathcal{T}$ . Now,

$$P(S(\mathbf{X}) = s) = \sum_{t \in \mathcal{T}} P(S(\mathbf{X}) = s | T(\mathbf{X}) = t) P_\theta(T(\mathbf{X}) = t).$$

Furthermore, since  $\sum_{t \in \mathcal{T}} P_\theta(T(\mathbf{X}) = t) = 1$ , we can write

$$P(S(\mathbf{X}) = s) = \sum_{t \in \mathcal{T}} P(S(\mathbf{X}) = s) P_\theta(T(\mathbf{X}) = t)$$

Therefore, if we define the statistic

$$g(t) = P(S(\mathbf{X}) = s | T(\mathbf{X}) = t) - P(S(\mathbf{X}) = s),$$

the above two equations show that

$$E_\theta g(T) = \sum_{t \in \mathcal{T}} g(t) P_\theta(T(\mathbf{X}) = t) = 0$$

for all  $\theta$ . Since  $T(\mathbf{X})$  is a complete statistic, this implies that  $g(t) = 0$  for all possible values  $t \in \mathcal{T}$ .

Hence, (6.2) is verified.  $\square$

It should be noted that the “minimality” of the sufficient statistic was not used in the proof of Basu’s theorem. Indeed, the theorem is true with this word omitted, because a fundamental property of a complete statistic is that it is minimal.

**Theorem 6.2.6 (Complete statistics in the exponential family)** *Let  $X_1, \dots, X_n$  be iid observations from an exponential family with pdf or pmf of the form*

$$f(\mathbf{x} | \boldsymbol{\theta}) = h(\mathbf{x}) c(\boldsymbol{\theta}) \exp\left(\sum_{j=1}^k w_j(\boldsymbol{\theta}) t_j(\mathbf{x})\right),$$

where  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_k)$ . Then the statistic

$$T(\mathbf{X}) = \left(\sum_{i=1}^n t_1(X_i), \sum_{i=1}^n t_2(X_i), \dots, \sum_{i=1}^n t_k(X_i)\right)$$

is complete if  $\{w_1(\boldsymbol{\theta}), \dots, w_k(\boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta\}$  contains an open set in  $\mathcal{R}^k$ .

**Example 6.2.15 (Using Basu's theorem)** Let  $X_1, \dots, X_n$  be iid exponential observations with parameter  $\theta$ . Consider computing the expected value of

$$g(\mathbf{X}) = \frac{X_n}{X_1 + \dots + X_n}.$$

We first note that the exponential distributions form a scale parameter family and thus  $g(\mathbf{X})$  is an ancillary statistic. The exponential distributions also form an exponential family with  $t(x) = x$  and so, by Theorem 6.2.6,

$$T(\mathbf{X}) = \sum_{i=1}^n X_i$$

is a complete statistic and by, Theorem 6.2.3,  $T(\mathbf{X})$  is a sufficient statistic. Hence, by Basu's theorem,  $T(\mathbf{X})$  and  $g(\mathbf{X})$  are independent. Thus we have

$$\theta = E_\theta X_n = E_\theta T(\mathbf{X})g(\mathbf{X}) = (E_\theta T(\mathbf{X}))(E_\theta g(\mathbf{X})) = n\theta E_\theta g(\mathbf{X}).$$

Hence, for any  $\theta$ ,  $E_\theta g(\mathbf{X}) = 1/n$ .

## 6.3 The likelihood Principle

### 6.3.1 The likelihood function

**Definition 6.3.1** Let  $f(\mathbf{x}|\theta)$  denote the joint pdf or pmf of the sample  $\mathbf{X} = (X_1, \dots, X_n)$ . Then, given that  $\mathbf{X} = \mathbf{x}$  is observed, the function of  $\theta$  defined by

$$L(\theta|\mathbf{x}) = f(\mathbf{x}|\theta)$$

is called the likelihood function.

This definition almost seems to be defining the likelihood function to be the same as the pdf or pmf. The only distinction between these two functions is which variable is considered fixed and which is varying. When we consider the pdf or pmf  $f(\mathbf{x}|\theta)$ , we are considering  $\theta$  as fixed and  $\mathbf{x}$  as the variable; when we consider the likelihood function  $L(\theta|\mathbf{x})$ , we are considering  $\mathbf{x}$  to be the observed sample point and  $\theta$  to be varying over all possible parameter values. If we compare the likelihood function at two parameter points and find that

$$L(\theta_1|\mathbf{x}) > L(\theta_2|\mathbf{x}),$$

then the sample we actually observed is more likely to have occurred if  $\theta = \theta_1$  than if  $\theta = \theta_2$ , which can be interpreted as saying that  $\theta_1$  is a more plausible value for the true value of  $\theta$  than is  $\theta_2$ . We carefully use the word “plausible” rather than “probable” because we often think of  $\theta$  as a fixed value. Furthermore, although  $f(\mathbf{x}|\theta)$ , as a function of  $\mathbf{x}$ , is a pdf, there is no guarantee that  $L(\theta|\mathbf{x})$ , as a function of  $\theta$ , is a pdf.

**LIKELIHOOD PRINCIPLE:** If  $\mathbf{x}$  and  $\mathbf{y}$  are two sample points such that  $L(\theta|\mathbf{x})$  is proportional to  $L(\theta|\mathbf{y})$ , that is, there exists a constant  $C(\mathbf{x}, \mathbf{y})$  such that

$$L(\theta|\mathbf{x}) = C(\mathbf{x}, \mathbf{y})L(\theta|\mathbf{y})$$

for all  $\theta$ , then the conclusions drawn from  $\mathbf{x}$  and  $\mathbf{y}$  should be identical.

Note that  $C(\mathbf{x}, \mathbf{y})$  may be different for different  $(\mathbf{x}, \mathbf{y})$  pairs but  $C(\mathbf{x}, \mathbf{y})$  does not depend on  $\theta$ .

## 6.4 The Equivariance Principle

The first type of equivariance might be called measurement equivariance. It prescribes that the inference made should not depend on the measurement scale that is used.

The second type of equivariance, actually an invariance, might be called formal invariance. It states that if two inference problems have the same formal structure in terms of the mathematical model used, then the same inference procedure should be used in both problems. The elements of the model that must be the same are:  $\Theta$ , the parameter space;  $\{f(\mathbf{x}|\theta) : \theta \in \Theta\}$ , the set of pdfs or pmfs for the sample; and the set of allowable inferences and consequences of wrong inferences.

For example,  $\Theta$  may be  $\Theta = \{\theta : \theta > 0\}$  in two problems. But in one problem  $\theta$  may be the average price of a dozen eggs and in another problem  $\theta$  may refer to the average height of giraffes. Yet, formal invariance equates these two parameter spaces since they both refer to the same set of real numbers.

**Equivariance Principle:** If  $Y = g(X)$  is a change of measurement scale such that the model for  $Y$  has the same formal structure as the model for  $X$ , then an inference procedure should be both measurement equivariant and formally equivariant.

**Example 6.4.1 (Binomial equivariance)** *Let  $X \sim \text{Binomial}(n, p)$  with known  $n$  and  $p$ . Let  $T(x)$  be the estimate of  $p$  that is used when  $X = x$  is observed. Rather than using the number of successes,  $X$ , to make an inference about  $p$ , we could use the number of failures,  $Y = n - X$ . We*

can see that  $Y \sim \text{Binomial}(n, q = 1 - p)$ . Let  $T^*(y)$  be the estimate of  $q$  that is used when  $Y = y$  is observed, so that  $1 - T^*(y)$  is the estimate of  $p$  when  $Y = y$  is observed. If  $x$  successes are observed, then the estimate of  $p$  is  $T(x)$ . But if there are  $x$  successes, then there are  $n - x$  failures and  $1 - T^*(n - x)$  is also an estimate of  $p$ . Measurement equivariance requires that these two estimates be equal, that is,  $T(x) = 1 - T^*(n - x)$ , since the change from  $X$  to  $Y$  is just a change in measurement scale. Furthermore, the formal structures of the inference problems based on  $X$  and  $Y$  are the same.  $X$  and  $Y$  both have  $\text{binomial}(n, \theta)$  distributions,  $0 \leq \theta \leq 1$ . So formal invariance requires that  $T(z) = T^*(z)$  for all  $z = 0, \dots, n$ . Thus, measurement and formal invariance together require that

$$T(x) = 1 - T^*(n - x) = 1 - T(n - x).$$

If we consider only estimators satisfying the above equality, then we have greatly reduced and simplified the set of estimators we are willing to consider. Whereas the specification of an arbitrary estimator requires the specification of  $T(0), \dots, T(n)$ , the specification of an estimator satisfying the above equality requires the specification only of  $T(0), \dots, T(\lfloor n/2 \rfloor)$ , where  $\lfloor n/2 \rfloor$  is the greatest integer not larger than  $n/2$ . This is the type of data reduction that is always achieved by the equivariance principle. The inference to be made for some sample points determines the inference to be made for other sample points.