

## Lecture 3: Sufficiency and Completeness

Instructor: Xiaohui Chen

Scribe: Yutong Dai

Last Modified: 2018-03-13

**Disclaimer:** *These notes have not been subjected to the usual scrutiny reserved for formal publications, and they may contain factual and/or typographic errors.*

**Motivation:** We want to summarize the information contained in data by determining a few key features of the sample values.

**Notation:** Denote  $X = (X_1, \dots, X_n)$ , where  $X_i \stackrel{iid}{\sim} f(x; \theta)$  and  $x$  denotes the realization of random samples.  $x$  also represents the variable of a function  $f(\cdot)$  and we will use  $X$  and  $x$  interchangeably.

Data reduction in terms of a particular statistic can be thought of a partition of the sample space  $\mathcal{X}$ . Let  $\mathcal{T} = \{t; t \in T(\mathcal{X})\}$ . Then  $T$  induces a partition of  $\mathcal{X}$ ,  $A_t$ , which is defined as  $A_t = \{x \in \mathcal{X}; T(x) = t, t \in \mathcal{T}\}$ . The statistic  $T(x)$  summarizes the data in that, rather than reporting the entire sample  $x$ , it reports only that  $T(x) = t$  or, equivalently,  $x \in A_t$ . For example, if  $T(x) = x_1 + \dots + x_n$ , then  $T(x)$  does not report the actual sample values but only the sum. There may be many different sample points that have the same sum.

We are interested in methods of data reduction that do not discard important information about the unknown parameter  $\theta$  and methods that successfully discard information that is irrelevant as far as gaining knowledge about  $\theta$  is concerned.

### 3.1 The Sufficiency Principle

A sufficient statistic for a parameter  $\theta$  is a statistic that, in a certain sense, captures all the information about  $\theta$  contained in the sample. Any additional information in the sample, besides the value of the sufficient statistic, does not contain any more information about  $\theta$ .

**Definition 3.1 (SUFFICIENCY PRINCIPLE)** If  $T(X)$  is a sufficient statistic(SS) for  $\theta$ , then any inference about  $\theta$  should depend on the sample  $X$  only through the value  $T(X)$ . That is, if  $x$  and  $y$  are two sample points such that  $T(x) = T(y)$ , then the inference about  $\theta$  should be the same whether  $X = x$  or  $Y = y$  is observed.

In this subsection we investigate some aspects of sufficient statistics and the Sufficiency Principle.

#### 3.1.1 Sufficient Statistics

**Definition 3.2 (SUFFICIENT STATISTICS)** A statistic  $T(X)$  is a sufficient statistic for  $\theta$  if the conditional distribution of the sample  $X$  given the value of  $T(X)$  does not depend on  $\theta$ .

To understand definition 3.2, let  $t$  be a possible value of  $T(X)$ . Of course we are interested in the conditional probability  $P_\theta(X = x | T(X) = t)$ . If the sample  $x$  is a point such that  $t \neq T(x)$ ,  $P_\theta(X = x | T(X) = t) = 0$ . Thus, we restrict  $t = T(x)$ , for some  $x \in \mathcal{X}$ . Then, by definition

$$P_\theta(X = x | T(X) = T(x))$$

does not depend on  $\theta$ , hence  $P_\theta(X = x|T(X) = T(x)) \triangleq P(X = x|T(X) = T(x))$ .

**Remark**

1. When  $T(X)$  has a continuous distribution, then  $P_\theta(T(X) = t) = 0$  for all values of  $t$ . A more sophisticated notion of conditional probability is needed. To simplify, we only focus on discrete case and point out analogous results that are true in the continuous case.
2. A statistic  $T(X)$  is not sufficient if  $P_\theta(X = x|T(X) = T(x))$  depends on  $\theta$ , for some  $x$ .

**Theorem 3.3** Suppose  $p(x; \theta)$  is the joint pmf of the sample  $X$  and  $q(t; \theta)$  is the pmf of  $T(X)$ . Then  $T(X)$  is a sufficient statistic for  $\theta \iff \frac{p(x; \theta)}{q(T(x); \theta)}$  is constant as a function of  $\theta$ .

**Proof:** " $\Rightarrow$ ":

$$P_\theta(X = x|T(X) = T(x)) = \frac{P_\theta(X = x, T(X) = T(x))}{P_\theta(T(X) = T(x))} \quad (3.1)$$

$$= \frac{P_\theta(X = x)}{P_\theta(T(X) = T(x))} \quad (3.2)$$

$$= \frac{p(x; \theta)}{q(T(x); \theta)} \quad (3.3)$$

where the second equality due to the fact that  $\{X; X = x\} \subset \{X; T(X) = T(x)\}$ . Since  $T(X)$  is a sufficient statistic for  $\theta$ ,  $\frac{p(x; \theta)}{q(T(x); \theta)}$  is constant as a function of  $\theta$ .

" $\Leftarrow$ " Just inverse the proof of the sufficiency. ■

**Remark**

1. It is still appropriate to use the above criterion to determine if  $T(X)$  is a sufficient statistic for  $\theta$  when  $X$  and  $T(X)$  have continuous distributions.
2. Sufficient statistics are not unique. Any one-to-one mappings of a sufficient statistic is also a sufficient statistic.

**Example 1**(Discrete case): Let  $X_i \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$ . Then the  $T(X) = \sum_{i=1}^n X_i$  is a sufficient statistic for  $\theta$ . Note that  $T(X)$  follows  $\text{Binomial}(n, \theta)$ . Then by Theorem 3.3, we know that

$$\frac{p(x; \theta)}{q(T(x); \theta)} = \frac{\prod_{i=1}^n \theta^{x_i} (1 - \theta)^{(1-x_i)}}{C_n^t \theta^t (1 - \theta)^{(n-t)}} \quad (t = \sum_i x_i) \quad (3.4)$$

$$= \frac{1}{C_n^t} \quad (3.5)$$

**Example 2**(Continuous case): Let  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$  and  $\sigma^2$  is known. Then the  $T(X) = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  is a sufficient statistic for  $\mu$ . Since  $\bar{X}_n$  follows  $N(\mu, \sigma^2/n)$ , we have

$$\frac{p(x; \mu)}{q(T(x); \mu)} = \frac{1}{\sqrt{n}} (2\pi\sigma^2)^{-(n-1)/2} \exp((n\bar{x}_n^2 - \sum_{i=1}^n x_i^2)/(2\sigma^2))$$

Then we show there exists a example that substantial data reduction is impossible.

**Example 3**(No reduction): Let  $X_i \stackrel{iid}{\sim} f(x)$ , where we can not specify more information about  $f$ . Then the ss can be  $T(X) = (X_{\pi(1)}, X_{\pi(2)}, \dots, X_{\pi(n)})$ , where  $\pi : \{1, 2, \dots, n\} \rightarrow \{1, 2, \dots, n\}$  is a permutation.

**Theorem 3.4 (Factorization Theorem)** Let  $f(x; \theta)$  denote the joint pdf or pmf of a sample  $X$ . A statistic  $T(X)$  is a sufficient for  $\theta \iff$  there exists functions  $g(t; \theta)$  and  $h(x)$  such that for all sample points  $x$  and all parameter points  $\theta$ ,

$$f(x; \theta) = h(x)g(T(x); \theta).$$

**Theorem 3.5 (SS for Exponential family)** Let  $X_1, \dots, X_n \stackrel{iid}{\sim} f(x; \theta)$  that belongs to an exponential family given by

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

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

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

is a sufficient statistic for  $\theta$ .

### 3.1.2 Minimal Sufficient Statistics

We have seen that SS are not unique, we introduce minimal sufficient statistics(MSS) to answer the questions that whether one SS is better than the other.

**Definition 3.6** A sufficient statistics  $M(X)$  is called a minimal sufficient statistic if for any other sufficient statistic  $T(X)$  there exists a functional  $g(\cdot)$  that  $M(X) = g(T(X))$ .

**Remark** We interpret the definition of MSS through the partitions on the sample space  $\mathcal{X}$ . Both  $T(X)$  and  $M(X)$  can induce partitions on  $\mathcal{X}$ , and we denote them as  $\{B_t = \{x \in \mathcal{X}; T(x) = t\}, t \in T(\mathcal{X})\}$  and  $\{A_m = \{x \in \mathcal{X}; M(x) = m\}, m \in M(\mathcal{X})\}$  receptively.  $\forall t \in T(\mathcal{X}), x \in B_t, \exists m \in M(\mathcal{X})$  such that  $g(T(x)) = m$ , hence  $B_t \subset A_m$ . That is any set in the partition induced by  $T$  is the subset of some set in the partition induced by  $M$ . Intuitively,  $M$  corresponds to the coarsest possible partition for  $\mathcal{X}$ , hence achieving the greatest possible data reduction for a SS.

**Theorem 3.7** Let  $f(x; \theta)$  be the pmf or the pdf of a sample  $X$ . Suppose there exists a function  $T(x)$  such that, for every two sample points  $x$  and  $y$ , the ratio  $\frac{f(x; \theta)}{f(y; \theta)}$  is constant as a function of  $\theta$  if and only if  $T(x) = T(y)$ . Then  $T(X)$  is a minimal sufficient statistic for  $\theta$ .

**Proof:** Without the loss of generality, we assume that  $f(x; \theta) > 0, \forall x \in \mathcal{X}, \theta \in \Theta$ .

Actually,  $T$  defines a partition on set  $\mathcal{X}$ . Let  $\mathcal{T} = \{t; t \in T(\mathcal{X})\}$ . Any fixed  $t \in \mathcal{T}$  corresponds to a set  $A_t = \{x \in \mathcal{X}; T(x) = t\}$ . Define  $g(t; \theta) = \sum_{x \in A_t} P_\theta(X = x) = \sum_{x \in A_t} f(x; \theta)$ . Now, we use factorization theorem to show that  $T(X)$  is a SS.

Note that

$$f(x; \theta) = g(T(x); \theta) \frac{f(x; \theta)}{g(T(x); \theta)} \quad (3.6)$$

$$= g(t; \theta) \frac{f(x; \theta)}{\sum_{y \in A_{T(x)}} f(y; \theta)} \quad (3.7)$$

$$= g(T(x); \theta) h(x) \quad (3.8)$$

Since  $h(x)$  is a constant in  $\theta$ ,  $T(X)$  is a SS. Suppose there exists another SS,  $T'(X)$ . By factorization theorem, we can find  $h'(X)$ ,  $g'(t; \theta)$  such that  $f(x; \theta) = g'(T'(x); \theta) h'(x)$ . Let  $x$  and  $y$  be any two sample points with  $T'(x) = T'(y)$ , then

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

Since  $T'$  is a SS, so  $h'$  does not depend on  $\theta$ , hence the ratio is a constant in  $\theta$ . By assumption, it implies that  $T(x) = T(y)$ . So  $T(X)$  is the MSS. ■

**Example 1:** Let  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$  and both  $\sigma^2, \mu$  are unknown. Let  $x$  and  $y$  denote two sample points and let  $(\bar{x}_n, S_x^2)$  and  $(\bar{y}_n, S_y^2)$  be the sample mean and variance respectively. Since

$$f(x; \mu, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2}[n(\bar{x}_n - \mu)^2 + \sum_i (x_i - \bar{x}_n)^2]\right).$$

Set  $g(t; \mu, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2}[n(t_1 - \mu)^2 + (n-1)t_2]\right)$  and  $h(x) = 1$ . We have  $f(x; \mu, \sigma^2) = g(T_1(x), T_2(x); \mu, \sigma^2) h(x)$ , where  $T_1(x) = \bar{x}_n$  and  $T_2(x) = S_x^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x}_n)^2$ . Therefore,

$$\frac{f(x; \mu, \sigma^2)}{f(y; \mu, \sigma^2)} = \exp\left(\frac{-n(\bar{x}_n - \bar{y}_n) + 2n\mu(\bar{x} - \bar{y}) - (n-1)(S_x^2 - S_y^2)}{2\sigma^2}\right).$$

This ratio will become constant in  $\mu$  and  $\sigma^2 \iff \bar{x} = \bar{y}$  and  $S_x^2 = S_y^2$ .

**Example 2:** Suppose  $X_i \stackrel{iid}{\sim} U[\theta, \theta + 1], \theta \in R$ . Then the joint distribution of a sample point  $x$  is

$$f(x; \theta) = \begin{cases} 1, & \theta \leq x_i \leq \theta + 1, \forall i \\ 0, & \text{otherwise} \end{cases}$$

$$\frac{f(x; \theta)}{f(y; \theta)} = \frac{\prod_{i=1}^n I(x_i \leq \theta + 1) I(x_i \geq \theta)}{\prod_{i=1}^n I(y_i \leq \theta + 1) I(y_i \geq \theta)} \text{ is constant in } \theta, \forall x, y \iff X_{(1)} = Y_{(1)} \text{ and } X_{(n)} = Y_{(n)}.$$

And thus  $T(X) = (X_{(1)}, X_{(n)})$ .

## 3.2 Ancillary Statistics

In the preceding sections, we consider sufficient statistics, which contain all the information about  $\theta$  that is available. Now, we consider another statistic - ancillary statistic.

**Definition 3.8** A statistic  $T(x)$  whose distribution doesn't depend on  $\theta$  is called an ancillary statistic.

From the Definition 3.8, we know that an ancillary statistic contains no information about  $\theta$ .

**Example 1:** Suppose  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$  and  $(\mu, \sigma^2)$  are both unknown. The statistic  $T(X) = \frac{X_1 - X_n}{\sqrt{S^2}}$  is an ancillary statistic, where  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ . This is because

$$T(X) = \frac{Z_1 - Z_2}{\sqrt{\sum_{i=1}^n (Z_i - \bar{Z}_n)^2}}$$

doesn't depend on  $\theta$ , where  $Z_i = \frac{X_i - \mu}{\sigma}$ . The location information is cancelled through taking difference and the scale information is eliminated by taking ratio.

**Example 2:** (Location-Scale Family) Suppose  $X_i \stackrel{iid}{\sim} \frac{1}{\sigma} f\left(\frac{x-\mu}{\sigma}\right)$  and  $(\mu, \sigma^2)$  are both unknown. The statistic  $T(X) = \left(\frac{X_1 - X_n}{\sqrt{S^2}}, \dots, \frac{X_{n-1} - X_n}{\sqrt{S^2}}\right)$  is an ancillary statistic.

However, following two examples will show that, even an ancillary statistic contains no information about the  $\theta$ , once we combine it with another statistic, which has some information about  $\theta$ , then we can have all information about  $\theta$ . Most surprisingly, we can use two ancillary statistic to derive a sufficient statistic.

**Example 3:** (Uniform ancillary)  $X_i \stackrel{iid}{\sim} U[\theta, \theta + 1], \theta \in R$ . Let  $X_{(1)}, \dots, X_{(n)}$  be the ordered statistic. Then  $R_n = X_{(n)} - X_{(1)}$  is the ancillary statistic.

**Example 4:**  $X_1, X_2 \stackrel{iid}{\sim} N(\theta, 1), \theta \in R$ .  $T_1 = X_1 - X_2$  is ancillary statistic. and  $T_2 = X_2$  isn't a sufficient statistic. But  $(T_1, T_2)$  is a sufficient statistic. Consider a one-to-one function  $f : (X_1 - X_2, X_2) \rightarrow (X_1, X_2)$ . Since  $X_1 - X_2$  is a sufficient statistic,  $(T_1, T_2)$  is also a sufficient statistic.

**Example 5:** Consider

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim N_2\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right).$$

Then  $T_1 = T_1(X, Y) = X, T_2 = T_2(X, Y) = Y$  are two ancillary statistic since their distributions don't depend on  $\rho$ . But  $(X, Y)$  is the a sufficient statistic for  $\rho$ .

### 3.3 Complete Statistic

On the one hand, a minimal sufficient statistic contains all information about the parameter  $\theta$ , hencing achieving the maximum data reduction. On the other hand, an ancillary statistci contains no information about the parameter  $\theta$ . So, intuitively, any minimal sufficient statistic should be functionally independent of any ancillary staistic. However, this claim doesn't necessary hold without any other constraint.

For example, consider  $X_i \stackrel{iid}{\sim} U[\theta, \theta + 1], \theta \in R$ . Then a minimal sufficient statistic is  $(X_{(n)}, X_{(1)})$ . Clearly, the  $(X_{(1)} - X_{(n)}, \frac{X_{(1)} + X_{(n)}}{2})$  is another minimal sufficient statistic. However, in this case  $X_{(n)} - X_{(1)}$  is ancillary and is an part of a minimal sufficiency statistic. Certainly, the anciallry and the minimal sufficient statistic are not functionally independent.

This example motivates us that we need put addiitonal requirements to ensure the independence.

**Definition 3.9 (Induced family of distributions)** Suppose  $X \sim f(x; \theta), \theta \in \Theta$  and  $T(X) \sim g(t; \theta)$ . The collection of pmf/pdf denoted by  $\{g(t; \theta); \theta \in \Theta\}$  is called the family of distribution induced by the statistic  $T(X)$ .

**Example:**  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ,  $\sigma^2$  is given and  $T(X) = \bar{X}_n = \frac{1}{n} \sum_i X_i \sim N(\mu, \sigma^2/n)$ . Then the induced family of distribution induced by  $T(X)$  is

$$\{g(t; \mu) = \frac{1}{\sqrt{2\pi\sigma^2/n}} \exp\left(\frac{-n(x - \mu)^2}{\sigma^2}\right); \mu \in R\}.$$

**Definition 3.10 (Completeness)** The family of distribution  $\{g(t; \theta); \theta \in \Theta\}$  induced by the statistic  $T(X)$  is complete if the following condition holds:

1.  $\forall$  real-valued function  $h(\cdot)$ , we have

$$E_\theta(|h(T)|) = \int |h(T(x))| f_X(x; \theta) dx < \infty \quad \forall \theta \in \Theta.$$

2. If  $\forall \theta \in \Theta, E_\theta(h(T)) = 0$  implies  $\forall \theta \in \Theta, P_\theta[h(T) = 0] = 1$ .

Moreover, a statistic  $T(X)$  is called complete iff the induced family of distribution induced by  $T(X)$  is complete.

**Example 1:** (Uniform complete) Suppose  $X_i \stackrel{iid}{\sim} U[0, \theta], \theta > 0$ .  $T(X) = \max_i X_i$ . Note that

$$F_T(t) = P(T(X) \leq t) = \begin{cases} 0 & , t < 0 \\ (t/\theta)^n & , 0 \leq t \leq \theta \\ 1 & , t > \theta, \end{cases}$$

we have

$$f_T(t) = \begin{cases} 0 & , t < 0 \\ n\theta^{-n}t^{n-1} & , 0 \leq t \leq \theta \\ 0 & , t > \theta. \end{cases}$$

Suppose  $\forall \theta \in \Theta$ ,

$$0 = E_\theta[h(T)] = \int_0^\theta h(t) n\theta^{-n} t^{n-1} dt,$$

which implies that  $\int_0^\theta h(t) t^{n-1} dt = 0, \forall \theta \in \Theta$ . Set derivative with respect to  $\theta$  on both side, we have  $h(\theta)\theta^{n-1} = 0, \forall \theta \in \Theta$ . This means that  $h(\theta) = 0, \forall \theta \in \Theta$ , and further more we have  $h(t) = 0, t \in (0, \theta], \forall \theta \in \Theta$  and  $\forall \theta \in \Theta, P_\theta[h(T) = 0] = 1$ .

**Example 2:** (Poisson complete) Let  $X_1, \dots, X_n$  denote a random sample from a Poisson distribution with parameter  $\lambda > 0$ . Then  $T(X) = \sum X_i \sim p(n\lambda)$ . For any function  $h(\cdot)$  and any  $\lambda$ ,

$$0 = E_\lambda(h(T(X))) = \sum_{t=0}^{+\infty} e^{-n\lambda} \frac{(n\lambda)^t}{t!} h(t).$$

This implies  $h(t) = 0, t = 0, 1, 2, \dots$  and therefore  $P_\lambda(h(T) = 0) = 1$ .

- In general a complete statistic is not necessary a sufficient statistic.

**Example 3:** (continued with **Example 2**) The joint distribution of  $(X_1, \dots, X_n)$  is

$$f(x_1, x_2, \dots, x_n; \lambda) = (e^{-n\lambda} \lambda^{\sum x_i}) \times \left( \frac{1}{x_1! x_2! \dots x_n!} \right) = g(T(x); \theta) h(x_1, x_2, \dots, x_n).$$

By the factorization theorem we know  $T(X)$  is a sufficient statistic. Now, consider  $T'(X) = X_1 + X_2$ . By the analysis in **Example 2**, we know that  $T'(X)$  is a complete statistic but not a sufficient statistic.

- However, a sufficient statistic is also not necessary a complete statistic.

**Example 4:** Let  $X_1, X_2$  denote a random sample from a Poisson distribution with parameter  $\lambda > 0$ . Then  $T(X) = (X_1, X_2)$  is sufficient but not complete. Consider a function  $h(x_1, x_2) = x_1 - x_2$ , then

$$E_\lambda(T(X)) = \sum_{x_1=0}^{+\infty} \sum_{x_2=0}^{+\infty} (x_1 - x_2) e^{-2\lambda} \frac{\lambda^{(x_1+x_2)}}{(x_1+x_2)!} = 0,$$

which implies  $P_\lambda(h(T) = 0) < 1$ . (Not that rigorous.)

Next, we introduce a theorem on the complete statistic for exponential family.

**Theorem 3.11 (exponential family)** Let  $X_1, \dots, X_n \stackrel{iid}{\sim} f(x; \theta) = h(x)c(\theta) \exp(\sum_{j=1}^k w_j(\theta)t_j(x))$ , where  $\theta = (\theta_1, \dots, \theta_k)^T$ . Then the statistic  $T(X) = (\sum_{i=1}^n t_1(X_i), \dots, \sum_{i=1}^n t_k(X_i))$  is a complete sufficient statistic, if the parameter space  $\{w_1(\theta), \dots, w_k(\theta); \theta \in \Theta\}$  contains a open subset.

**Remark:**

1. From the Theorem 3.11, we know that for curved exponential family, the sufficient statistic is typically not complete.
2. Since the transformation  $T(\cdot)$  is usually continues, (not that rigorous) the condition that the parameter space  $\{w_1(\theta), \dots, w_k(\theta); \theta \in \Theta\}$  contains a open subset is equivalent to the one that  $\Theta$  contains a open subset.

**Example :** (Simple Linear Regression) Consider the following linear regression setting:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

where  $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$ ,  $\sigma^2$  is known and  $\theta = (\beta_0, \beta_1)$ . The joint pdf of  $(y_1, \dots, y_n)$  is

$$f(y_1, \dots, y_n; \theta) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{\sum_{i=1}^n (\beta_0 + \beta_1 x_i)^2}{2\sigma^2}\right) \exp\left(-\frac{\sum_{i=1}^n y_i^2}{2\sigma^2}\right) \exp\left(\frac{2 \sum_{i=1}^n (\beta_0 + \beta_1 x_i) y_i}{2\sigma^2}\right) \quad (3.9)$$

$$= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{\sum_{i=1}^n (\beta_0 + \beta_1 x_i)^2}{2\sigma^2}\right) \exp\left(-\frac{\sum_{i=1}^n y_i^2}{2\sigma^2}\right) \exp\left(\frac{\beta_0}{\sigma^2} \sum_{i=1}^n y_i\right) \exp\left(\frac{\beta_1}{\sigma^2} \sum_{i=1}^n x_i y_i\right) \quad (3.10)$$

$$= h(x)c(\theta) \exp(w_1(\theta)t_1(y)) \exp(w_2(\theta)t_2(y)). \quad (3.11)$$

Since  $\{(\frac{\beta_0}{\sigma^2}, \frac{\beta_1}{\sigma^2}); \beta_0, \beta_1 \in R\}$  contains a open subset, we know  $T(Y) = (\sum_{i=1}^n Y_i, \sum_{i=1}^n x_i Y_i)$  is a complete sufficient statistic for  $(\beta_0, \beta_1)$ .

**Remark:** Recall that the LS estimator for  $(\beta_0, \beta_1)$  is

$$\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \begin{pmatrix} n & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum_i Y_i \\ \sum_i x_i Y_i \end{pmatrix} = AT(Y).$$

Since  $A$  is invertible, which means  $A$  is a one-to-one transformation, we conclude that  $(\hat{\beta}_0, \hat{\beta}_1)$  is a complete sufficient statistic estimator. This implies that LS estimator is really good in the sense of completeness and sufficiency.

Now, with the help of completeness, Basu gives the following theorem to show that when the sufficient statistic is independent of ancillary statistic.

**Theorem 3.12 (Basu's Theorem)** If  $T(X)$  is a complete and sufficient statistic, then  $T(X)$  is independent of every ancillary statistic.

**Proof:** We only give the proof of discrete case.

Let  $S(X)$  be any ancillary statistic. We want to prove that  $\forall s, t$ ,

$$P_\theta(S(X) = s | T(X) = t) = P_\theta(S(X) = s) = P(S(X) = s),$$

where the last equality holds as  $S(X)$  is ancillary statistic. Since the  $T(X)$  is the sufficient, we know that

$$P_\theta(S(X) = s | T(X) = t) = P(S(X) = s | T(X) = t).$$

So the only thing left is to prove that  $\forall s, t$

$$P(S(X) = s | T(X) = t) - P(S(X) = s) = 0.$$

Define a function

$$g_s(t) = P(S(X) = s | T(X) = t) - P(S(X) = s),$$

which doesn't depend on  $\theta$ . Take expectation with respect to  $T(X)$ , we have

$$E_\theta[g_s(T)] = \sum_t g_s(t) P_\theta(T(X) = t) \tag{3.12}$$

$$= \sum_t P(S(X) = s | T(X) = t) P_\theta(T(X) = t) - \sum_t P(S(X) = s) P_\theta(T(X) = t) \tag{3.13}$$

$$= P(S(X) = s) - P(S(X) = s) \tag{3.14}$$

$$= 0 \tag{3.15}$$

By the completeness of  $T(X)$ , we know that  $P_\theta(g_s(T) = 0) = 1, \forall \theta$ , which completes the proof. ■

**Example :**  $X_i \stackrel{iid}{\sim} N(0, \sigma^2), \sigma^2 > 0$  but is unknown. Let  $S^2 = \frac{(\sum_{i=1}^n x_i)^2}{\sum_{i=1}^n x_i^2}$ . What is the value of  $E(S^2)$  ?

Consider  $T(X) = \sum_{i=1}^n x_i^2$ , then it's easy to see that  $T(X)$  is complete and sufficient. Since  $S^2$  belongs to the scale family, hence it is ancillary. By Theorem 3.12, we know that  $T(X) \perp S(X)$ . Then  $E(T(X))E(S^2) = E[(\sum_{i=1}^n x_i^2)]$ , which implies that  $E(S^2) = 1$ .



## 3.4 Appendix

### 3.4.1 pdfs of order statistic

Suppose the cdf of  $X$  is  $F_X(x)$  and  $X_{(1)}, \dots, X_{(n)}$  are the ordered statistic, then

$$f_{X_{(k)}}(x) = \frac{n!}{(k-1)!(n-k)!} [F_X(x)]^{k-1} [1 - F_X(x)]^{n-k} f_X(x) \quad (3.16)$$

$$f_{X_{(j)}, X_{(k)}}(x, y) = \frac{n!}{(j-1)!(k-j-1)!(n-k)!} \times [F_X(x)]^{j-1} [F_X(y) - F_X(x)]^{k-j-1} [1 - F_X(y)]^{n-k} f_X(x) f_X(y) \text{ where } x \leq y \quad (3.17)$$

$$f_{X_{(1)}, \dots, X_{(n)}}(x_1, \dots, x_n) = n! f_X(x_1) \cdots f_X(x_n) \text{ where } x_1 \leq x_2 \leq \cdots \leq x_n. \quad (3.18)$$