

# Sufficient Statistics

A sample is a sequence  $X_1, \dots, X_n$  of iid random variables.

Convention : We use capital letters to denote the sample as random variable, and lower case letters for the observations. Thus the sample  $X_1, \dots, X_n$  consists of random variables while  $x_1, \dots, x_n$  denotes the observed (measured) values of the sample.

Assume the common distribution of  $X_1, X_2, \dots, X_n$  depends on a parameter  $\theta$ . Here  $\theta$  can be one parameter or multiple parameters.

Def A statistics (or estimator of  $\theta$ ) is a function of the sample space:

$$T = u(X_1, \dots, X_n)$$

Def We call  $T$  a sufficient statistics if the conditional distribution of the sample given  $T$  does not depend on the parameter  $\theta$ .

In other words, if  $f(x; \theta)$  is the distribution of  $X$  (either pmf or pdf), then

$$\frac{f(x_1; \theta) \dots f(x_n; \theta)}{f_T(u(x_1, \dots, x_n); \theta)} = \underbrace{h(x_1, \dots, x_n)}_{\text{independent of } \theta}$$

Alternatively we can formulate this as (2)

$(X_1, \dots, X_n) | T=t$  does not depend on  $\theta$ .  
This, from the definition of the conditional expectation is

$$\frac{f_{X_1}(x_1; \theta) \dots f_{X_n}(x_n; \theta)}{f_T(t; \theta)} \text{ is independent of } \theta \text{ if } T(x_1, \dots, x_n) = t.$$

This is a canonical way of defining a sufficient statistics.

An alternative way of describing this is the following theorem.

Th (Neyman)  $T = u(X_1, \dots, X_n)$  is a sufficient statistics if and only if

$$f(x_1; \theta) \dots f(x_n; \theta) = k_1(u(x_1, \dots, x_n); \theta) k_2(x_1, \dots, x_n)$$

where  $k_2$  is a function which depends only on  $x_1, \dots, x_n$  and not on  $\theta$ .

Pf We will do a proof in the case of discrete i.v.

The direct implication is straightforward because

$$f(x_1; \theta) \dots f(x_n; \theta) = f_T(u(x_1, \dots, x_n); \theta) H(x_1, \dots, x_n)$$

For the reverse if we have

$$f(x_1; \theta) \dots f(x_n; \theta) = k_1(u(x_1, \dots, x_n); \theta) k_2(x_1, \dots, x_n)$$

with  $u(x_1, \dots, x_n) = t$ , we have that

$$f_{X|T}(x_1, \dots, x_n) = \frac{f(x_1; \theta) \dots f(x_n; \theta)}{f_T(t; \theta)} = \frac{k_1(t; \theta)}{k_2(x_1, \dots, x_n)}$$

This means that for fixed  $x_1, \dots, x_n$  such that  $x_1 + \dots + x_n = t$ ,  $\frac{k_1(t; \theta)}{f_T(t; \theta)}$  does not depend on  $\theta$ . Thus

$$\frac{k_1(t; \theta)}{f_T(t; \theta)} = g(t) \text{ and thus}$$

$$f(x_1; \theta) \dots f(x_n; \theta) = f_T(t; \theta) k_3(x_1, \dots, x_n)$$

Def We say that an estimator  $T$  is unbiased if  $E[T] = \theta$

The estimator is called MUE if for any other unbiased estimator  $T'$  we have

$$\text{Var}(T) \leq \text{Var}(T')$$

Given a sufficient statistics and an unbiased estimator  $T$ , then we can find another one with a smaller variance.

Th ( Rao-Blackwell ) If  $T_1$  is a sufficient statistics and  $T_2$  is an unbiased est, then  $T_3 = E[T_2 | T_1]$  is actually also unbiased and has smaller variance than  $T_1$ .

Pf We know that  $E[T_3] = E[T_2] = \theta$ .  $T_3$  is unbiased. On the other hand,

(See the Appendix here)

$$\text{Var}(T_2) = \text{Var}(T_3) + E[(T_2 - T_3)^2] \geq \text{Var}(T_3)$$

Equality is attained iff  $T_3 = T_2$  or in other words if  $T_2$  is a function of  $T_1$ . If we prevent this, then  $T_3$  has a strictly smaller variance.

## Completeness and Uniqueness

~~A~~ A variable  $Z$  from  $(f(x; \theta))_\theta$  is called complete if any function  $u$  such that  $E_\theta[u(Z)] = 0 \forall \theta$  implies  $u \equiv 0$  a.s..

Ex Take  $Z \sim \text{Poisson}(\theta)$ . Then  $f(x; \theta) = e^{-\theta} \frac{\theta^x}{x!}$  is a complete family

Sol For a function  $u$ , we have

$$\begin{aligned} E[u(Z)] &= \sum_{x=0}^{\infty} u(x) P(Z=x) \\ &= \sum_{x=0}^{\infty} u(x) \frac{\theta^x}{x!} e^{-\theta} = 0 \quad \forall \theta \end{aligned}$$

Therefore  $\sum_{x=0}^{\infty} u(x) \frac{\theta^x}{x!} = 0 \quad \forall \theta > 0$ , so  $u \equiv 0$

### Theorem (Lehmann-Scheffé)

If  $X_1, \dots, X_n$  is a sample and a sufficient statistic  $Y = u(X_1, \dots, X_n)$  has a family  $(f_Y(y; \theta))$  which is complete, and there is a function  $Y_2 = \varphi(Y_1)$  which is an unbiased estimator of  $\theta$ , then  $Y_2$  is the unique MVUE.

Pf From Rao-Blackwell we know that we can restrict our search to  $Y_3 = \varphi(Y_1)$ .

If  $Y_3$  is also unbiased, then  $E[\varphi(Y_1)] = \theta = E[\varphi(Y_1)]$  and thus  $E[\varphi(Y_1) - \varphi(Y_1)] = 0$ . From completeness we get that  $\varphi = \varphi$  a.s., thus  $Y_2 = Y_3$ .

Slang We say that  $Y_1$  is a complete sufficient statistic.

## Appendix Conditioning

Typically if we have  $X$  and  $Y$ , two random variables then the distribution of  $X$  given  $Y$  is defined via the conditional pmf or pdf.

$$P_{X|Y}(x|y) = \frac{P(X=x, Y=y)}{P(Y=y)} = \frac{P_{X,Y}(x,y)}{P_Y(y)}$$

Inspired by this, if  $(X,Y)$  have a joint density, then

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}.$$

This is good if we do have a joint density.

Ex What would be the conditional dist of  $X$  given  $X$ ?

$$P_{X|X}(x|y) = \frac{P(X=x, X=y)}{P(X=y)} = \begin{cases} 0 & \text{if } x \neq y \\ 1 & \text{if } x = y \end{cases}$$

What this means is that the distribution of  $X$  given  $X=y$  is actually concentrated at  $y$ .

This is the case in the discrete situation, though in the continuous case is even more problematic because  $f_{X|X}(x,y)$  does not really make sense.

However the conditional expectation of any two random variables actually makes sense.

More precisely if we have  $X$  and  $Y$ , we can try to define  $E[X|Y]$  as the best approximation of  $X$  in terms of  $Y$ . In other words we want  $E[X|Y] = \varphi(Y)$  such that

$E[(X - \varphi(Y))^2]$  has the smallest possible value among all possible functions  $\varphi$ .

This leads to some properties of  $\varphi(Y)$ .

Properties ①  $E[(X - \varphi(Y))\psi(Y)] = 0 \quad \forall \psi$   
 In particular  $E[X] = E[\varphi(Y)]$   
 ②  $\text{Var}(\varphi(Y)) \leq \text{Var}(X)$

Pf ① If we assume a  $\varphi$  which minimizes  $E[(X - \varphi(Y))^2]$  over all choices of  $\varphi$  then by replacing  $\varphi$  by  $\varphi + \varepsilon\psi$  and taking derivatives w.r.t  $\varepsilon$  at  $\varepsilon = 0$  we get

$$E[(X - \varphi(Y))\psi(Y)] = 0$$

Alternatively we can consider the Hilbert space of all functions  $\varphi(Y)$  such that  $E[\varphi(Y)^2] < \infty$

Then one instance of  $\varphi(Y)$  is obtained by projecting  $X$  from  $L^2(\Omega, \mathcal{F}, P)$  onto this subspace.

$$\begin{aligned} \textcircled{2} \text{Var}(X) &= E[(X - E[X])^2] = \\ &= E[(X - \varphi(Y) + \varphi(Y) - E[X])^2] \\ &= E[(X - \varphi(Y))^2] + 2 \underbrace{E[(X - \varphi(Y))(\varphi(Y) - E[X])]}_{\substack{\text{from the first} \\ \text{part}}} + E[(\varphi(Y) - E[X])^2] \end{aligned}$$

Because  $E\{x\} = E\{\varphi(Y)\}$  we actually get

$$\text{Var}(x) = E[(x - \varphi(Y))^2] + \text{Var}(\varphi(Y)).$$

In particular we actually get something better, namely  $\text{Var}(x) \geq \text{Var}(\varphi(Y))$  with equality iff  $x = \varphi(Y)$  which means that  $x$  is a function of  $Y$ .

Def The  $\varphi(Y)$  defined above is denoted  $E\{x|Y\}$ .

The function  $\varphi$  itself is denoted as

$$E\{x|Y=y\} = \varphi(y)$$

Ex Compute  $E\{x|x^2\}$  if  $x \sim N(0,1)$ .

Sol  $E\{x|x^2\} = \varphi(x^2)$  and to find this  $\varphi$  we can use the first property above, namely

$$E\{x \psi(x^2)\} = E\{\varphi(x^2) \psi(x^2)\} \text{ for any } \psi.$$

In particular

$$\underbrace{\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x \psi(x^2) e^{-x^2/2} dx}_0 = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \varphi(x^2) \psi(x^2) e^{-x^2/2} dx$$

thus  $E\{x|x^2\} = 0$ .

Interestingly,  $E\{x^2|x\} = x^2$  because  $x^2$  is already a function of  $x$ .

Using this interpretation of the conditional expectation we actually have also that a formal definition can be put forward in the form:

Def  $E[X|Y] = Z$  such that

$$E[Z \psi(Y)] = E[X \psi(Y)]$$

for any  $\psi$  continuous and bounded.

Th If  $X$  is integrable, then  $E[X|Y]$  always exists.

In the case  $X, Y$  have a joint pdf then

$$Z = \varphi(Y) \text{ with}$$

$$\varphi(y) = \frac{\int x f_{X,Y}(x,y) dx}{f_Y(y)}$$

because in this case

$$\begin{aligned} E[\varphi(Y) \psi(Y)] &= \iint \frac{x f_{X,Y}(x,y)}{f_Y(y)} \psi(y) f_Y(y) dy \\ &= \iint x \psi(y) f_{X,Y}(x,y) dx dy \\ &= E[X \psi(Y)]. \end{aligned}$$

thus we fall back onto the well known case.



Exercise A rigorous way of formulating that

$(X_1, \dots, X_n) / U(X_1, \dots, X_n)$  does not depend on  $\theta$  is to state that

$E[\xi(X_1, \dots, X_n) | U(X_1, \dots, X_n)] \stackrel{= \varphi_3(U(X_1, \dots, X_n))}{\text{does not}} \text{ on } \theta \text{ for any choice of } X_1, \dots, X_n.$

By definition, this is the same as saying that

$$E[\xi(X_1, \dots, X_n) \psi(U(X_1, \dots, X_n))] = E[\underbrace{\varphi_3(U(X_1, \dots, X_n))}_{\xi} \psi(U(X_1, \dots, X_n))] \text{ for any choice of } \psi.$$

The right hand side is the same as

$$\int \xi(x_1, \dots, x_n) \psi(U(x_1, \dots, x_n)) f(x_1; \theta) \dots f(x_n; \theta) dx_1 \dots dx_n \\ = \int \varphi_3(t) \psi(t; \theta) f_T(t; \theta) dt$$

Taking  $\psi(t; \theta) = \frac{\psi(t)}{f_T(t; \theta)}$  we get that

$$\int \xi(x_1, \dots, x_n) \psi(U(x_1, \dots, x_n)) \frac{f(x_1; \theta) \dots f(x_n; \theta)}{f_T(U(x_1, \dots, x_n); \theta)} dx_1 \dots dx_n \\ = \int \varphi_3(t) \psi(t) dt$$

Since  $\xi$  is an arbitrary function of  $x_1, \dots, x_n$  and  $\varphi_3, \psi$  do not depend on  $\theta$ , it must be that

$$\frac{f(x_1; \theta) \dots f(x_n; \theta)}{f_T(U(x_1, \dots, x_n); \theta)} \text{ does not depend on } \theta.$$