

Sample: $X_1, \dots, X_n \stackrel{iid}{\sim} X$ with pdf $f(x|\theta)$

$\tilde{X} = (X_1, \dots, X_n)^T$ n-dimensional random vector.

$\underline{x} = (x_1, \dots, x_n)^T$ n-dimensional real-valued vector.

Ex. $n=2$. $X_1, X_2 \sim N(\mu, \sigma^2)$

$\tilde{X} = (X_1, X_2)^T$

observed $(10.5, 10.9)$ $\underline{x} = (10.5, 10.9)^T$

Statistics: $T(\tilde{X})$: a form of data reduction or data summary

Sample Space \longrightarrow Image of T

\mathcal{X} $\longrightarrow \mathcal{Y} = \{t : t = T(\underline{x}) \text{ for some } \underline{x} \in \mathcal{X}\}$

Example. $n=2$. $N(\mu, \sigma^2)$

$\mathcal{X} = \mathbb{R}^2$

$T(\underline{x}) = \frac{x_1 + x_2}{2}$ $\mathbb{R}^2 \rightarrow \mathbb{R}$

Image of T : \mathbb{R}

Goal: find $T(\tilde{X})$
for θ

Sufficiency Principle.. If $T(\underline{X})$ is a sufficient statistic for θ ,
 then any inference about θ should depend on the sample \underline{X}
 only through the value $T(\underline{X})$

if we have observed two samples \underline{x} and \underline{y}
 such that $T(\underline{x}) = T(\underline{y})$, then the inference about θ
 should be the same whether \underline{x} or \underline{y} is observed!

Eg: $N(\theta, 1)$ point estimator: \bar{X}_n
 confidence interval $\bar{X}_n \pm z_{\alpha/2} \frac{1}{\sqrt{n}}$

Definition of Sufficient Statistics

A statistic $T(\underline{X})$ is a sufficient statistic for θ
 if the conditional distribution of the sample \underline{X} given
 the value of $T(\underline{X})$ does not depend on θ .

$$\underline{X} = (X_1, \dots, X_n)^\top \quad \prod_{i=1}^n f(x_i | \theta)$$

How to find a sufficient statistic for θ ?

Theorem 6.2.2 if $p(\underline{x}|\theta)$ is the joint pdf or pmf of \underline{X}
and $q(t|\theta)$ is the pdf or pmf of $T(\underline{X})$

then $T(\underline{X})$ is a sufficient statistic for θ

if the ratio

$$\frac{p(\underline{x}|\theta)}{q(t(\underline{x})|\theta)}$$
 is free of θ .

Example: (Binomial Sufficient Statistic)

Let $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta) \quad 0 < \theta < 1$

$$T(\underline{X}) = \sum_{i=1}^n X_i \sim \text{Binomial}(n, \theta)$$

$$p(\underline{x}|\theta) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} = \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} = \theta^{\sum x_i} (1-\theta)^{n-\sum x_i}$$

$$q(t(\underline{x})|\theta) = \binom{n}{t(\underline{x})} \theta^{t(\underline{x})} (1-\theta)^{n-t(\underline{x})} \quad t(\underline{x}) = \sum x_i$$

$$= \binom{n}{\sum x_i} \theta^{\sum x_i} (1-\theta)^{n-\sum x_i}$$

$$\chi = \left\{ (x_1, \dots, x_n) : x_i = 1 \text{ or } 0 \right\}$$

$$\frac{P(\mathbf{x}|\theta)}{q(t(\mathbf{x})|\theta)} = \frac{\theta^{\sum x_i} (1-\theta)^{n-\sum x_i}}{\binom{n}{\sum x_i} \theta^{\sum x_i} (1-\theta)^{n-\sum x_i}} \quad \text{for all } \mathbf{x} \in \mathcal{X}$$

$$= \frac{1}{\binom{n}{\sum x_i}}$$

free of θ

Thus $T(\mathbf{x}) = \sum x_i$ is a sufficient statistic for θ

Example: $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, \sigma^2)$ where σ^2 is known
what is a sufficient statistic for θ ?

Guess: try $T(\mathbf{x}) = \bar{X}_n \sim N(\theta, \frac{\sigma^2}{n})$

$$P(\mathbf{x}|\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_i - \theta)^2}{2\sigma^2}\right\} \quad \mathbf{x} \in \mathbb{R}^n$$

$$q(t(\mathbf{x})|\theta) = \frac{1}{\sqrt{2\pi\sigma^2/n}} \exp\left\{-\frac{(\frac{1}{n}\sum x_i - \theta)^2}{2\sigma^2/n}\right\} \quad \mathbf{x} \in \mathbb{R}^n$$

$$\begin{aligned} \frac{P(\mathbf{x}|\theta)}{q(t(\mathbf{x})|\theta)} &\propto \exp\left(-\frac{\sum (x_i - \theta)^2}{2\sigma^2}\right) \\ &\quad \exp\left(-\frac{(\frac{1}{n}\sum x_i - \theta)^2}{2\sigma^2/n}\right) \\ &= \exp\left(-\left[\sum_{i=1}^n (x_i - \frac{1}{n}\sum x_i)^2 + n(\frac{1}{n}\sum x_i - \theta)^2\right] / (2\sigma^2)\right) \\ &= \exp\left(-n(\frac{1}{n}\sum x_i - \theta)^2 / (2\sigma^2)\right) \end{aligned}$$

$$= \exp\left(-\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{2\sigma^2}\right)$$

free of θ !!!

Ex (Sufficient order statistic)

$$X_1, \dots, X_n \sim f(x|\theta)$$

$$p(\underline{x}|\theta) = \prod_{i=1}^n f(x_i|\theta)$$

$$\underline{T(X)} = (X_{(1)}, \dots, X_{(n)})^\top \quad \text{order statistic}$$

$$q(T(X)|\theta) = n! \prod_{i=1}^n f(x_{(i)}|\theta) = n! \prod_{i=1}^n f(x_i|\theta)$$

p.

ratio

$$\frac{1}{n!} \quad \text{free of } \theta !$$

Theorem 6.2.6 (Factorization Theorem)

Let $p(\underline{x}|\theta)$ be the joint pdf/pmf of \underline{X}

a statistic $T(\underline{X})$ is a sufficient statistic for θ

if and only if

there exist function $g(T(\underline{x})|\theta)$ and $h(\underline{x})$

such that

$$p(\underline{x}|\theta) = g(T(\underline{x})|\theta) h(\underline{x}) \quad \text{for all } \underline{x} \in \mathcal{X}$$

Example (Uniform Sufficient statistic)

$$X_1, \dots, X_n \stackrel{iid}{\sim} \text{Unif}(1, \theta)$$

$$\begin{aligned} p(\mathbf{x}|\theta) &= \prod_{i=1}^n \frac{1}{\theta-1} \mathbb{I}(1 < X_i < \theta) \\ &= \left(\frac{1}{\theta-1}\right)^n \prod_{i=1}^n \mathbb{I}(1 < X_i < \theta) \\ &= \left(\frac{1}{\theta-1}\right)^n \mathbb{I}(1 < X_{(1)} < X_{(n)} < \theta) \\ &= \underbrace{\left(\frac{1}{\theta-1}\right)^n \mathbb{I}(X_{(n)} < \theta)}_{g(t(\mathbf{x})|\theta)} \times \underbrace{\mathbb{I}(1 < X_{(1)})}_{h(\mathbf{x})} \end{aligned}$$

$$t(\mathbf{x}) = X_{(n)}$$

Eg $X_1, \dots, X_n \sim \text{Unif}(\theta_1, \theta_2) \quad \theta = \left(\begin{matrix} \theta_1 \\ \theta_2 \end{matrix}\right)$

$$p(\mathbf{x}|\theta) = \underbrace{\left(\frac{1}{\theta_2 - \theta_1}\right)^n \mathbb{I}(\theta_1 < X_{(1)} < X_{(n)} < \theta_2)}_{g(t(\mathbf{x})|\theta)} \times 1$$

$$\downarrow \quad g(t(\mathbf{x})|\theta) \times h(\mathbf{x})$$

$$t(\mathbf{x}) = \begin{pmatrix} X_{(1)} \\ X_{(n)} \end{pmatrix}$$