

Math 4501 - Probability and Statistics II

6.7 - Sufficient Statistics

Relation with maximum likelihood estimators

Theorem

Let X_1, X_2, \dots, X_n denote a random sample from a distribution with pmf/pdf $f(x; \theta)$, with $\theta \in \Omega$. \leftarrow unknown parameter:

If a sufficient statistic $Y = u(X_1, X_2, \dots, X_n)$ for θ exists and if a maximum likelihood estimator $\hat{\theta}$ of θ also exists uniquely, then $\hat{\theta}$ is a function of $Y = u(X_1, X_2, \dots, X_n)$.

REMARK: $\hat{\theta} = v(Y)$ \leftarrow this theorem + v is invertible \rightarrow $\hat{\theta}$ is also sufficient

INTERPRETATION: MLE are based on sufficient statistic
and (under mild conditions)
MLE are sufficient too!

Proof sketch: x_1, \dots, x_m random sample from distribution $f(x; \theta)$, $\theta \in \Omega$

$Y = u(x_1, \dots, x_m)$ is sufficient iff

FISHER-NEYMAN FACTORIZATION THM

$$f_{\text{joint}}(\underline{x_1, \dots, x_m}; \theta) = \underbrace{\phi(y, \theta)}_{\phi \text{ depends on } x_1, \dots, x_m \text{ only through}} \cdot \underbrace{h(x_1, \dots, x_m)}_{h \text{ does not depend on } \theta}$$

BUT, then the likelihood function is

$$L(\theta) = \phi(y, \theta) \cdot \underbrace{h(x_1, \dots, x_m)}_{\text{does not depend on } \theta}$$

Thus, to maximize $L(\theta)$ w.r.t θ , we need to maximize $\phi(y, \theta)$ w.r.t θ
 When we maximize ϕ w.r.t θ , we obtain $\hat{\theta} = \hat{\theta}(y)$, which by assumption is unique \Rightarrow MLE is a function of the sufficient statistic!!

More than one unknown parameter

The notions and results discussed above admit a natural extension to the case of two or more unknown parameters.

We illustrate such extension by discussing the case of two unknown parameters.

Definition (Jointly sufficient statistics)

Let X_1, \dots, X_n be a random sample from a distribution with pdf/pmf $f(\cdot; \theta_1, \theta_2)$ depending on two parameters $(\theta_1, \theta_2) \in \Omega$.

The statistics $Y_1 = u_1(X_1, \dots, X_n)$ and $Y_2 = u_2(X_1, \dots, X_n)$ are said to be jointly sufficient for θ_1 and θ_2 if the conditional distribution of X_1, \dots, X_n given $Y_1 = y_1$ and $Y_2 = y_2$ does not depend on (θ_1, θ_2) for any values y_1 of Y_1 and y_2 of Y_2 .

Same interpretation as in the one-parameter case: once the sufficient statistics are given, there is no additional information about the parameters left in the remaining (conditional) distribution.

$$P(X_1 = x_1, \dots, X_n = x_n \mid Y_1 = y_1, Y_2 = y_2) \text{ does not depend on } \theta_1 \text{ and } \theta_2$$

Theorem (Fisher-Neyman Factorization Theorem: two parameters)

Let X_1, X_2, \dots, X_n denote random variables with joint pdf/pmf $f(x_1, x_2, \dots, x_n; \theta_1, \theta_2)$ depending on the parameters θ_1 and θ_2 .

The statistics $Y_1 = u_1(X_1, X_2, \dots, X_n)$ and $Y_2 = u_2(X_1, X_2, \dots, X_n)$ are jointly sufficient for θ_1 and θ_2 if and only if

$$f(x_1, \dots, x_n; \theta_1, \theta_2) = \phi(u_1(x_1, \dots, x_n), u_2(x_1, \dots, x_n); \theta_1, \theta_2) \cdot h(x_1, \dots, x_n),$$

does not depend on θ_1 or θ_2

where ϕ depends on x_1, \dots, x_n only through $u_1(x_1, \dots, x_n)$ and $u_2(x_1, \dots, x_n)$ while $h(x_1, \dots, x_n)$ does not depend upon θ_1 or θ_2 .

Note: as in the one-parameter case, it is often easier to check sufficiency using the Factorization Theorem than it is using the definition.

Remark

Let $\underline{Y_1}$ and $\underline{Y_2}$ be jointly sufficient for θ_1 and θ_2 , and let

$$(Z_1, Z_2) = \nu(\underline{Y_1}, \underline{Y_2})$$

be in invatible

be a function of $\underline{Y_1}$ and $\underline{Y_2}$ not involving θ_1 and θ_2 that has a single-valued inverse.

Then $\underline{Z_1}$ and $\underline{Z_2}$ are also jointly sufficient statistics for θ_1 and θ_2 .

Example

Let X_1, X_2, \dots, X_n be a random sample from $N(\theta_1, \theta_2)$, $\theta_1 \in \mathbb{R}$, $\theta_2 > 0$.

Show

$$Y_1 = \sum_{i=1}^n X_i \quad \text{and}$$

$$Y_2 = \sum_{i=1}^n X_i^2$$

pdf of $N(\theta_1, \theta_2)$ is

$$f(x; \theta_1, \theta_2) = \frac{1}{\sqrt{2\pi\theta_2}} e^{-\frac{(x-\theta_1)^2}{2\theta_2}}$$

are jointly sufficient statistics for θ_1 and θ_2 .

Conclude that

$$\bar{X} \text{ is MLE for } \mu \rightarrow \bar{X} = \frac{Y_1}{n} \quad \text{and} \quad S^2 = \frac{Y_2 - \frac{1}{n} Y_1^2}{n-1} \quad \leftarrow s^2 \text{ is a function of MLE for } \sigma^2$$

are also jointly sufficient statistics for θ_1 and θ_2

To show that Y_1 and Y_2 are jointly sufficient for θ_1 and θ_2 , we will use the Fisher-Neyman factorization theorem:

$$f_{\text{joint}}(x_1, \dots, x_n; \theta_1, \theta_2) = \prod_{i=1}^n f(x_i; \theta_1, \theta_2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\theta_2}} e^{-\frac{(x_i-\theta_1)^2}{2\theta_2}}$$

$$= \left(\frac{1}{\sqrt{2\pi\theta_2}} \right)^m \exp \left(- \sum_{i=1}^m \frac{(x_i - \theta_1)^2}{2\theta_2} \right) =$$

$$= (2\pi\theta_2)^{-m/2} \exp \left(- \frac{1}{2\theta_2} \sum_{i=1}^m (x_i - \theta_1)^2 \right)$$

$$= (2\pi\theta_2)^{-m/2} \exp \left(- \frac{1}{2\theta_2} \sum_{i=1}^m (x_i^2 - 2\theta_1 x_i + \theta_1^2) \right)$$

$$= (2\pi\theta_2)^{-m/2} \exp \left(- \frac{1}{2\theta_2} \sum_{i=1}^m x_i^2 + \frac{\theta_1}{\theta_2} \sum_{i=1}^m x_i - \frac{m\theta_1^2}{2\theta_2} \right)$$

$$= (2\pi\theta_2)^{-m/2} \exp \left(\frac{\theta_1}{\theta_2} \sum_{i=1}^m x_i - \frac{1}{2\theta_2} \sum_{i=1}^m x_i^2 - \frac{m\theta_1^2}{2\theta_2} \right)$$

$\hookrightarrow \phi(Y_1, Y_2, \theta_1, \theta_2)$

$\cdot \underbrace{1}_{h(x_1, \dots, x_n)}$

CONCLUSION:

$$f_{\text{joint}}(x_1, \dots, x_m; \theta_1, \theta_2) = f(y_1, y_2, \theta_1, \theta_2) \cdot h(x_1, \dots, x_m)$$

with $f(y_1, y_2, \theta_1, \theta_2) = (2\pi\theta_2)^{-m/2} \cdot \exp\left(\frac{\theta_1}{\theta_2}y_1 - \frac{1}{2\theta_2}y_2 - \frac{m\theta_1^2}{2\theta_2}\right)$

$$y_1 = \underbrace{\sum_{i=1}^m x_i}, \quad y_2 = \sum_{i=1}^m x_i^2$$

$$h(x_1, \dots, x_m) = 1$$

which means that

$$y_1 = \sum_{i=1}^m x_i \quad \text{and} \quad y_2 = \sum_{i=1}^m x_i^2$$

are jointly sufficient for θ_1 and θ_2

To see that \bar{X} and s^2 are also sufficient, it is enough to check that \bar{X} and s^2 may be written as an invertible function of y_1 and y_2

that is, there exists an invertible function φ such that

$$(\bar{X}, s^2) = \varphi(y_1, y_2)$$

Note that:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} y_1$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i^2 - 2x_i \bar{X} + \bar{X}^2)$$

$$= \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - 2\bar{X} \underbrace{\sum_{i=1}^n x_i}_{m\bar{X}} + m\bar{X}^2 \right) = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - m\bar{X}^2 \right)$$

$$= \frac{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2}{n-1}$$

y_1

CONCLUSION:

$$\bar{X} = \frac{1}{m} Y_1$$

$$\rightarrow S^2 = \frac{Y_2 - \frac{1}{m} (Y_1)^2}{m-1}$$

inverse in

$$Y_1 = m \bar{X}$$

$$Y_2 = (m-1) S^2 + m (\bar{X})^2$$

$\Rightarrow \bar{X}$ and S^2 are
jointly sufficient
statistics for μ and σ^2

is invertible

because we can solve the system

$$x_1 = \frac{1}{m} y_1$$

$$x_2 = \frac{y_2 - y_1^2/m}{m-1}$$

for y_1 and y_2 in terms x_1 and x_2
and the solution is unique

$$\text{that is } (\bar{X}, S^2) = \mathcal{V}(Y_1, Y_2)$$

$$\text{where } \mathcal{V}(Y_1, Y_2) = \left(\frac{1}{m} Y_1, \frac{Y_2 - Y_1^2/m}{m-1} \right)$$

Math 4501 - Probability and Statistics II

7.1 - Confidence intervals for means

Overview

We will see how to determine interval estimates for the unknown mean of a certain distribution.

We will require that the actual unknown mean falls within such intervals with a specified likelihood, representing the level of confidence in our interval estimate:

- we call such interval estimates confidence intervals.

exact results { We will construct confidence intervals for the mean of:

- a normal distribution with known variance
- a normal distribution with unknown variance

approximation { We will construct approximate confidence intervals for the mean of:

- an unknown distribution, using a large sample (known and unknown variance)
- an unknown distribution, using a smaller sample ↵ smaller but not necessarily small!!

Normal distribution with known variance

Let X_1, X_2, \dots, X_n be a random sample from a normal distribution $N(\mu, \sigma^2)$ with known variance σ^2 .

\bar{X} is the MLE for μ

unknown
known!!

We will use \bar{X} , the unbiased estimator of the distribution mean μ , to determine an interval estimate for μ , called a *confidence interval* for μ .

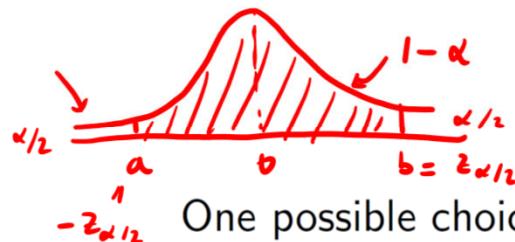
We have previously seen that \bar{X} is $N(\mu, \sigma^2/n)$ and so $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ is $N(0, 1)$.

$$x_1, \dots, x_n \sim N(\mu, \sigma^2)$$

$$\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$$

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

Take a small number $\alpha \in (0, 1)$, and look for two numbers a and b such that



One possible choice is to take $a = -z_{\alpha/2}$ and $b = z_{\alpha/2}$ to get

$$P\left(a \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq b\right) = 1 - \alpha$$

90%
95%
98%
99%

$$P\left(-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2}\right) = 1 - \alpha$$

→ solve
inequality
inside for μ !

Solve the inequality in

for μ to get

$$P \left(-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2} \right) = 1 - \alpha$$

$$-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2}$$

multiply by $\frac{\sigma}{\sqrt{n}}$

$$-z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right) \leq \bar{X} - \mu \leq z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)$$

add $-\bar{X}$ to all terms

$$-\bar{X} - z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right) \leq -\mu \leq -\bar{X} + z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)$$

multiply by -1
(invert inequality)

$$\bar{X} - z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right) \leq \mu \leq \bar{X} + z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)$$

Since the probability of the first inequality is $1 - \alpha$, so is the probability of the last.

Thus, we have

$$P\left[\bar{X} - z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}}\right) \leq \mu \leq \bar{X} + z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}}\right)\right] = 1 - \alpha,$$

and we conclude that the probability of the random interval

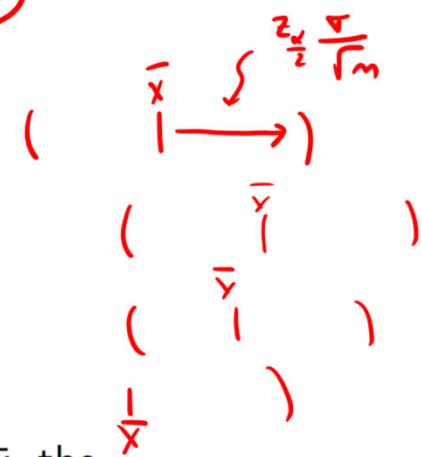
$$\left[\bar{X} - z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}}\right), \bar{X} + z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}}\right)\right]$$

including the unknown mean μ is $1 - \alpha$

Once the sample is observed and the sample mean computed to equal \bar{x} , the interval

$$\left[\bar{x} - z_{\alpha/2} \left(\sigma / \sqrt{n}\right), \bar{x} + z_{\alpha/2} \left(\sigma / \sqrt{n}\right)\right]$$

becomes known.



Since the probability that the random interval

$$[\bar{x} - z_{\alpha/2}(\sigma/\sqrt{n}), \bar{x} + z_{\alpha/2}(\sigma/\sqrt{n})]$$

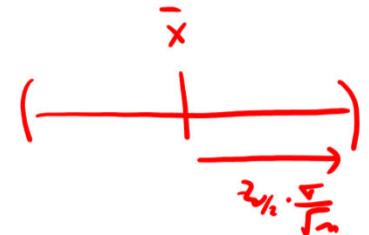
covers μ before the sample is drawn is equal to $1 - \alpha$, we call the computed interval a $100(1 - \alpha)\%$ confidence interval for the unknown mean μ

The number $100(1 - \alpha)\%$, or just, $1 - \alpha$, is called the confidence coefficient.

Notes:

1) The confidence interval for μ :

$$[\bar{x} - z_{\alpha/2}(\sigma/\sqrt{n}), \bar{x} + z_{\alpha/2}(\sigma/\sqrt{n})]$$



is centered at the point estimate \bar{x} and has length $2z_{\alpha/2}(\sigma/\sqrt{n})$.

2) As n increases, $z_{\alpha/2}(\sigma/\sqrt{n})$ decreases, resulting in a shorter confidence interval with the same confidence coefficient $1 - \alpha$.

- A shorter confidence interval gives a more precise estimate of μ , regardless of the confidence we have in the estimate of μ .
- There is a potential cost (time, money, effort) in increasing the sample size, even when more observations may be collected.

3) For a fixed sample size n , the length of the confidence interval can also be shortened by decreasing the confidence coefficient $1 - \alpha$.

- If this is done, we achieve a shorter confidence interval at the expense of losing some confidence.

Example

Let \underline{X} equal the length of life of a 60-watt light bulb marketed by a certain manufacturer.

normal distribution with known variance

Assume that the distribution of X is $N(\mu, 1296)$. A random sample of $n = 27$ bulbs is tested until they burn out, yielding a sample mean of $\bar{x} = 1478$ hours.

Determine a 95% confidence interval for μ .

Key reasoning now to set up CI:

Note that:

- (1) • sample is taken from a normal distribution
- (2) • variance is known $\sigma^2 = 1296 \Rightarrow \sigma = \sqrt{1296} = 36$

The 95% confidence interval for μ is

$$\begin{aligned} & \left[\bar{x} - z_{0.025} \left(\frac{\sigma}{\sqrt{n}} \right), \bar{x} + z_{0.025} \left(\frac{\sigma}{\sqrt{n}} \right) \right] = \\ &= \left[1478 - 1.96 \left(\frac{36}{\sqrt{27}} \right), 1478 + 1.96 \left(\frac{36}{\sqrt{27}} \right) \right] = [1464.42, 1491.58]. \end{aligned}$$

$$\begin{aligned} & \text{---} \quad \sigma^2 \text{ is known} \\ & \Rightarrow X_1, \dots, X_n \sim N(\mu, \sigma^2) \Rightarrow \bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \\ & \Rightarrow Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \quad \alpha/2 \quad 1 - \alpha \quad \alpha/2 \\ & \Rightarrow P\left(-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2}\right) = 1 - \alpha \end{aligned}$$

Endpoints are
 $\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$