

# Statistical Methods for Data Science

Lessons 17 and 18 - Confidence intervals: normal data, large sample method, linear regression.

Salvatore Ruggieri

Department of Computer Science  
University of Pisa  
[salvatore.ruggieri@unipi.it](mailto:salvatore.ruggieri@unipi.it)

# From point estimate to interval estimate

- For an estimator  $T = h(X_1, \dots, X_n)$  and a dataset  $x_1, \dots, x_n$ , we derive a *point estimate*:

$$t = h(x_1, \dots, x_n)$$

- Sometimes, a *range* of plausible values for an unknown parameter is preferred
- Idea: *confidence interval* is an interval for which we can be confident the unknown parameter is in with a specified probability (*confidence level*)

# Example

- From the Chebyshev's inequality:

$$P(|Y - \mu| < k\sigma) \geq 1 - \frac{1}{k^2}$$

For  $Y = \bar{X}_n$ ,  $k = 2$  and  $\sigma = 100$  Km/s:

$$P(|\bar{X}_n - \mu| < 200) \geq 0.75$$

- ▶ i.e.,  $\bar{X}_n \in (\mu - 200, \mu + 200)$  with probability  $\geq 75\%$  [random variable in a fixed interval]
  - ▶ or,  $\mu \in (\bar{X}_n - 200, \bar{X}_n + 200)$  with probability  $\geq 75\%$  [fixed value in a random interval]
- $(\bar{X}_n - 200, \bar{X}_n + 200)$  is an interval estimator of the unknown  $\mu$
- Let  $t = 299852.4$  be the estimate (realization of  $T = \bar{X}_n$ )
- $\mu \in (t - 200, t + 200) = (299652.4, 300052.4)$  is correct with confidence  $\geq 75\%$

Table 17.1. Michelson data on the speed of light.

850	740	900	1070	930	850	950	980	980	880
1000	980	930	650	760	810	1000	1000	960	960
960	940	960	940	880	800	850	880	900	840
830	790	810	880	880	830	800	790	760	800
880	880	880	860	720	720	620	860	970	950
880	910	850	870	840	840	850	840	840	840
890	810	810	820	800	770	760	740	750	760
910	920	890	860	880	720	840	850	850	780
890	840	780	810	760	810	790	810	820	850
870	870	810	740	810	940	950	800	810	870

# More info on $T$ means better (smaller) intervals

- Assume  $X_i \sim N(\mu, \sigma^2)$ . Hence,  $\bar{X}_n \sim N(\mu, \sigma^2/n)$  and:

$$Z_n = \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \sim N(0, 1)$$

- $P(-1.15 \leq Z_n \leq 1.15) = \Phi(1.15) - \Phi(-1.15) = 0.75$ 
  - ▶  $-1.15 = q_{0.125}$  and  $1.15 = q_{0.875}$  are called the critical values for achieving 75% probability
- Going back to  $\bar{X}_n$ :

$$P\left(\bar{X}_n - 1.15 \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + 1.15 \frac{\sigma}{\sqrt{n}}\right) = 0.75$$

- $\mu \in \left(t - 1.15 \frac{200}{\sqrt{100}}, t + 1.15 \frac{200}{\sqrt{100}}\right) = \mathbf{(t-23, t+23)}$  is correct *with confidence = 75%*

# Confidence intervals

CONFIDENCE INTERVALS. Suppose a dataset  $x_1, \dots, x_n$  is given, modeled as realization of random variables  $X_1, \dots, X_n$ . Let  $\theta$  be the parameter of interest, and  $\gamma$  a number between 0 and 1. If there exist sample statistics  $L_n = g(X_1, \dots, X_n)$  and  $U_n = h(X_1, \dots, X_n)$  such that

$$P(L_n < \theta < U_n) = \gamma$$

for every value of  $\theta$ , then

$$(l_n, u_n),$$

where  $l_n = g(x_1, \dots, x_n)$  and  $u_n = h(x_1, \dots, x_n)$ , is called a  $100\gamma\%$  confidence interval for  $\theta$ . The number  $\gamma$  is called the *confidence level*.

- Sometimes, only have  $P(L_n < \theta < U_n) \geq \gamma$  [*conservative*  $100\gamma\%$  confidence interval]
  - ▶ E.g., the interval found using Chebyshev's inequality
- There is no way of knowing if  $l_n < \theta < u_n$  (interval is correct)
- We only know that we have probability  $\gamma$  of covering  $\theta$
- Notation:  $\gamma = 1 - \alpha$  where  $\alpha$  is called the *significance level*
  - ▶ E.g., 95% confidence level equivalent to 0.05 significance level

**Seeing theory simulation**

# Confidence interval for the mean

- Let  $X_1, \dots, X_n$  be a random sample and  $\mu = E[X_i]$  to be estimated
- Problem: confidence intervals for  $\mu$  ?
  - ▶ Normal data
    - with known variance
    - with unknown variance
  - ▶ General data (with unknown variance)
    - large sample, i.e., large  $n$
    - bootstrap (next lesson)

# Critical values

## Critical value

The (right) *critical value*  $z_p$  of  $Z \sim N(0, 1)$  is the number with right tail probability  $p$ :

$$P(Z \geq z_p) = p$$

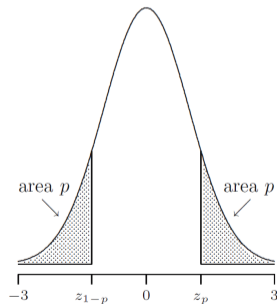
- Alternatively,  $p = 1 - \Phi(z_p) = 1 - P(Z \leq z_p)$ .
  - ▶ This is why Table B.1 of the textbook is given for  $1 - \Phi()$
- Alternatively,  $\Phi(z_p) = 1 - p$ , i.e.,  $z_p$  is the  $(1 - p)$ th quantile
- By  $P(Z \geq z_p) = P(Z \leq -z_p) = p$ , we have:

$$P(Z \geq -z_p) = 1 - P(Z \leq -z_p) = 1 - p$$

and then:

$$z_{1-p} = -z_p$$

- ▶ E.g.,  $z_{0.975} = -z_{0.025} = -1.96$  and  $z_{0.025} = -z_{0.975} = 1.96$



# CI for the mean: normal data with known variance

- Dataset  $x_1, \dots, x_n$  realization of random sample  $X_1, \dots, X_i \sim N(\mu, \sigma^2)$
- Estimator  $\bar{X}_n \sim N(\mu, \sigma^2/n)$  and the scaled mean:

$$Z = \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \sim N(0, 1) \quad (1)$$

- Confidence interval for  $Z$ :

$$P(c_l \leq Z \leq c_u) = \gamma \quad \text{or} \quad P(Z \leq c_l) + P(Z \geq c_u) = \alpha = 1 - \gamma$$

- Symmetric split:

$$P(Z \leq c_l) = P(Z \geq c_u) = \alpha/2$$

Hence  $c_u = -c_l = z_{\alpha/2}$ , and by (1):

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

$(\bar{X}_n - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X}_n + z_{\alpha/2} \frac{\sigma}{\sqrt{n}})$  is a  $100\gamma\%$  or  $100(1 - \alpha)\%$  confidence interval for  $\mu$



# One-sided confidence intervals

- One-sided confidence intervals (*greater-than*):

$$P(L_n < \theta) = \gamma$$

Then  $(l_n, \infty)$  is a  $100\gamma\%$  or  $100(1 - \alpha)\%$  one-sided confidence interval

- $l_n$  is called the *lower confidence bound*
- Normal data with known variance:

$$P(\bar{X}_n - z_\alpha \frac{\sigma}{\sqrt{n}} \leq \mu) = 1 - \alpha = \gamma$$

$(\bar{X}_n - z_\alpha \frac{\sigma}{\sqrt{n}}, \infty)$  is a  $100\gamma\%$  or  $100(1 - \alpha)\%$  one-sided confidence interval for  $\mu$

**See R script**

# CI for the mean: normal data with unknown variance

- When  $\sigma^2$  is **unknown**, we use its unbiased estimator:  $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$
- The following transformation is called *studentized mean*:

$$T = \sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \sim t(n-1)$$

DEFINITION. A continuous random variable has a *t-distribution with parameter*  $m$ , where  $m \geq 1$  is an integer, if its probability density is given by

$$f(x) = k_m \left(1 + \frac{x^2}{m}\right)^{-\frac{m+1}{2}} \quad \text{for } -\infty < x < \infty,$$

where  $k_m = \Gamma\left(\frac{m+1}{2}\right) / \left(\Gamma\left(\frac{m}{2}\right) \sqrt{m\pi}\right)$ . This distribution is denoted by  $t(m)$  and is referred to as the *t-distribution with  $m$  degrees of freedom*.

- Student t-distribution  $X \sim t(m)$ :
  - ▶ for  $m = 1$ , it is the Cauchy distribution
  - ▶  $E[X] = 0$  for  $m \geq 2$ , and  $\text{Var}(X) = m/(m-2)$  for  $m \geq 3$
  - ▶ For  $m \rightarrow \infty$ ,  $X \rightarrow N(0, 1)$

See R script

# CI for the mean: normal data with unknown variance

- Dataset  $x_1, \dots, x_n$  realization of random sample  $X_1, \dots, X_i \sim N(\mu, \sigma^2)$

## Critical value

The (right) *critical value*  $t_{m,p}$  of  $T \sim t(m)$  is the number with right tail probability  $p$ :

$$P(T \geq t_{m,p}) = p$$

- Same properties of  $z_p$
- From the studentized mean:

$$T = \sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \sim t(n-1)$$

to confidence interval:

$$P\left(\bar{X}_n - t_{n-1, \alpha/2} \frac{S_n}{\sqrt{n}} \leq \mu \leq \bar{X}_n + t_{n-1, \alpha/2} \frac{S_n}{\sqrt{n}}\right) = 1 - \alpha = \gamma$$

$(\bar{X}_n - t_{n-1, \alpha/2} \frac{S_n}{\sqrt{n}}, \bar{X}_n + t_{n-1, \alpha/2} \frac{S_n}{\sqrt{n}})$  is a  $100\gamma\%$  or  $100(1 - \alpha)\%$  confidence interval for  $\mu$

**See R script**

# CI for the mean: general data with unknown variance

- Dataset  $x_1, \dots, x_n$  realization of random sample  $X_1, \dots, X_i \sim N(\mu, \sigma^2)$
- A variant of CLT states that for  $n \rightarrow \infty$

$$T = \sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \rightarrow N(0, 1)$$

- For large  $n$ , we make the approximation:

*[how large should  $n$  be?]*

$$T = \sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \approx N(0, 1)$$

and then

$$P\left(\bar{X}_n - z_{\alpha/2} \frac{S_n}{\sqrt{n}} \leq \mu \leq \bar{X}_n + z_{\alpha/2} \frac{S_n}{\sqrt{n}}\right) \approx 1 - \alpha = \gamma$$

$(\bar{X}_n - z_{\alpha/2} \frac{S_n}{\sqrt{n}}, \bar{X}_n + z_{\alpha/2} \frac{S_n}{\sqrt{n}})$  is a  $100\gamma\%$  or  $100(1 - \alpha)\%$  confidence interval for  $\mu$

**See R script**

# Determining the sample size

- The narrower the CI the better (smaller variability)
- Sometimes, we start with an accuracy requirement:
  - ▶ find a  $100(1 - \alpha)\%$  CI  $(l_n, u_n)$  such that  $u_n - l_n \leq w$
- How to set  $n$  to satisfy the  $w$  bound?
- Case: normal data with known variance
  - ▶ CI is  $(\bar{X}_n - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X}_n + z_{\alpha/2} \frac{\sigma}{\sqrt{n}})$
  - ▶ Bound on the CI is:

$$2z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq w$$

leading to:

$$n \geq \left(2z_{\alpha/2} \frac{\sigma}{w}\right)^2$$

# CI for other parameters: linear regression coefficients

- Simple linear regression:  $Y_i = \alpha + \beta x_i + U_i$  with  $U_i \sim \mathcal{N}(0, \sigma^2)$
- We have  $\hat{\beta} \sim \mathcal{N}(\beta, \text{Var}(\hat{\beta}))$  where  $\text{Var}(\hat{\beta}) = \sigma^2 / SXX$  is unknown
- The studentized statistics is  $t(n-2)$ -distributed:

*[prove it!]*  
*[proof omitted]*

$$\frac{\hat{\beta} - \beta}{\sqrt{\text{Var}(\hat{\beta})}} \sim t(n-2)$$

- For  $\gamma = 0.95$ :

$$P(-t_{n-2,0.025} \leq \frac{\hat{\beta} - \beta}{\sqrt{\text{Var}(\hat{\beta})}} \leq t_{n-2,0.025}) = 0.95$$

and then a 95% confidence interval is:  $\hat{\beta} \pm t_{n-2,0.025} \text{se}(\hat{\beta})$  where  $\text{se}(\hat{\beta}) = \hat{\sigma}^2 / \sqrt{SXX}$

- Similarly, we get for  $\alpha$ ,  $\hat{\alpha} \pm t_{n-2,0.025} \text{se}(\hat{\alpha})$

**See R script**

## CI for other parameters: expectation of fitted values

- Simple linear regression:  $Y_i = \alpha + \beta x_i + U_i$  with  $U_i \sim \mathcal{N}(0, \sigma^2)$
- For the fitted values  $\hat{y} = \hat{\alpha} + \hat{\beta}x_0$  at  $x_0$ , a 95% confidence interval is:

$$\hat{y} \pm t_{n-2,0.025} se(\hat{Y})$$

where  $se(\hat{Y}) = \hat{\sigma} \sqrt{\left(\frac{1}{n} + \frac{(\bar{x}_n - x_0)^2}{SXX}\right)}$

- This interval concerns *the expectation of fitted values* at  $x_0$ .
  - ▶ E.g., the mean of predicted values at  $x_0$  is in  $[\hat{y} + t_{n-2,0.025} se(\hat{Y}), \hat{y} - t_{n-2,0.025} se(\hat{Y})]$

**See R script**

## CI for other parameters: fitted values

- Simple linear regression:  $Y_i = \alpha + \beta x_i + U_i$  with  $U_i \sim \mathcal{N}(0, \sigma^2)$
- For a given single prediction, we must also account for the error term  $U$  in:

$$\hat{V} = \hat{\alpha} + \hat{\beta}x_0 + U$$

- Assuming  $U \sim \mathcal{N}(0, \sigma^2)$ , we have  $\text{Var}(\hat{V}) = \sigma^2(1 + \frac{1}{n} + \frac{(\bar{x}_n - x_0)^2}{SXX})$
- A 95% confidence interval is:

$$\hat{y} \pm t_{n-2, 0.025} \text{se}(\hat{V})$$

$$\text{where } \text{se}(\hat{V}) = \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(\bar{x}_n - x_0)^2}{SXX}}$$

- A predicted value at  $x_0$  is in  $[\hat{y} + t_{n-2, 0.025} \text{se}(\hat{V}) \text{ and } \hat{y} - t_{n-2, 0.025} \text{se}(\hat{V})]$

**See R script**