# Confidence intervals for linear regression parameters:

A ==confidence interval (CI)== at the $c = 100(1-\alpha)\%$ level refers to an interval around a stochastic quantity in which it will lie $c\%$ of times. A related quantity is the significance level which is denoted as $\alpha$.

| Confidence level $c$ | Significance level $\alpha$ |
|---|---|
| 99% | 0.01 |
| 95% | 0.05 |
| 90% | 0.1 |

Hypothesis testing refers to the validation of a certain statement

statement called the alternate hypothesis (Ha), with respect to a given significance level.

Variable of interest $\longrightarrow$ Population mean of the linear regression parameters.

Hypothesis tests for the population mean
- Z-test (population $\sigma$ is known)
- t-test (==population $\sigma$ is unknown==)

## Z-test

Let us consider an independent sample $(y_1, y_2, \cdots, y_n)$ collected from a population with a known population standard deviation

Sample mean $\quad \bar{y} = \dfrac{1}{n} \displaystyle\sum_{i=1}^{n} \boxed{y_i}$

Sample variance $\quad \sigma^2 = \dfrac{1}{n-1} \displaystyle\sum_{i=1}^{n} (y_i - \bar{y})^2$

$\rightarrow E\left[\bar{y}\right] = E\left[\dfrac{1}{n}\displaystyle\sum_{i=1}^{n} y_i\right] = \dfrac{1}{n}\displaystyle\sum_{i=1}^{n} E[y_i] = \dfrac{1}{\cancel{n}} \times \cancel{n} \times \mu = \boxed{\mu}$

$\rightarrow Var\left[\bar{y}\right] = Var\left[\dfrac{1}{n}\displaystyle\sum_{i=1}^{n} y_i\right] = \dfrac{1}{n^2}\displaystyle\sum_{i=1}^{n} Var[y_i] = \dfrac{1}{n^2} \times n\sigma^2 = \boxed{\dfrac{\sigma^2}{n}}$
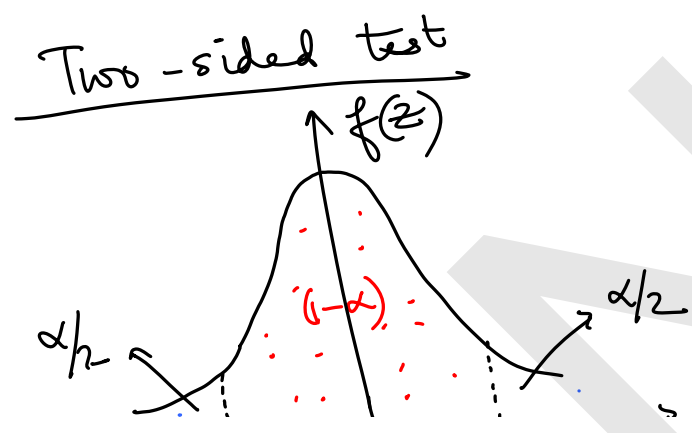
If we assume each $y_i$ to be normally distributed, then:

$$y_i \sim N(\mu, \sigma^2)$$

$$\bar{y} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

$$Z = \frac{\bar{y} - \mu}{\sigma/\sqrt{n}} \implies \boxed{Z \sim \mathcal{N}(0, 1)}$$

Z-statistic

Hypothesis test $\longrightarrow$ One-sided test

$\longrightarrow$ Two-sided test

Two-sided test

$f(z)$

$(1-\alpha)$

$\alpha/2$

$\alpha/2$

$H_0$ (null hypothesis) is rejected if

$$Z > Z_{1-\frac{\alpha}{2}} \quad \text{or} \quad Z < Z_{\frac{\alpha}{2}}$$

Otherwise the null hypothesis cannot be

## One-sided test



$H_0$ (null hypothesis) is rejected if

$$Z > Z_{1-\alpha}$$

Otherwise the null hypothesis cannot be rejected.

## Student's t test

Test for the population mean when the population standard deviation ($\sigma$) is unknown.

Random variable

$$T = \frac{\bar{y} - \mu}{s/\sqrt{n}} \sim t_{n-1}$$

sample standard deviation

Student's t distribution with $(n-1)$ degrees of freedom.

$$P\left( t_{\frac{\alpha}{2}} \leq T \leq t_{1-\frac{\alpha}{2}} \right) = 1-\alpha$$

$$P\left( t_{\frac{\alpha}{2}} \leq \left( \frac{\overline{y}-\mu}{s/\sqrt{n}} \right) \leq t_{1-\frac{\alpha}{2}} \right) = 1-\alpha$$

$$P\left( t_{\frac{\alpha}{2}} \frac{s}{\sqrt{n}} \leq \overline{y} - \mu \leq t_{1-\frac{\alpha}{2}} \cdot \frac{s}{\sqrt{n}} \right) = 1-\alpha$$

$$P\left( -t_{1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}} \leq \mu - \overline{y} \leq -t_{\frac{\alpha}{2}} \frac{s}{\sqrt{n}} \right) = 1-\alpha$$

$$P\left( \overline{y} - t_{1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}} \leq \mu \leq \overline{y} - t_{\frac{\alpha}{2}} \frac{s}{\sqrt{n}} \right) = 1-\alpha$$
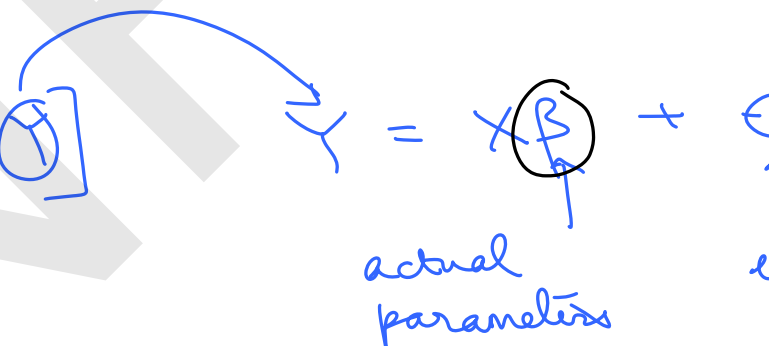
# Confidence intervals [CI] for linear regression parameters:

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

To assign CI to $\hat{\beta}$, we need to know:

- $E(\hat{\beta})$
- $Var(\hat{\beta})$

We will assume that the errors $(\epsilon)$ in prediction are normally distributed.

$$E(\hat{\beta}) = E\left[(X^T X)^{-1} X^T Y\right]$$

$$Y = X \beta + \epsilon$$

actual parameters      error

$$E\left[\hat{\beta}\right] = E\left[(X^TX)^{-1}X^T(X\beta + \epsilon)\right]$$

$$= E\left[\underbrace{(X^TX)^{-1}X^TX}_{I}\beta + \underline{(X^TX)^{-1}X^T\epsilon}\right]$$

$$= E\left[I\beta\right] + (X^TX)^{-1}X^T \underbrace{E(\epsilon)}_{\to 0}$$

$$= E\left[\beta\right] + 0$$

$$E\left[\hat{\beta}\right] = E\left[\beta\right] = \beta$$

$\hat{\beta}$ is an <u>unbiased estimator</u> for $\beta$.

$$Var\left(\hat{\beta}\right) = ?$$

↳ To determine this, we will use the covariance matrix

$$\text{Cov}(\hat{\beta}) = E\left[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^T\right]$$

$$Y = X\beta + \epsilon$$

$$\hat{\beta} = (X^TX)^{-1}X^TY \implies X^TX\hat{\beta} = \boxed{X^TY}$$

$$\boxed{X^TY} = X^TX\beta + X^T\epsilon$$

$$X^TX\hat{\beta} = X^TX\beta + X^T\epsilon$$

$$\implies \hat{\beta} - \beta = (X^TX)^{-1}X^T\epsilon$$

$$\text{Cov}(\hat{\beta}) = E\left[(X^TX)^{-1}X^T\epsilon \cdot \left((X^TX)^{-1}X^T\epsilon\right)^T\right]$$

$$= (X^TX)^{-1} E\left[\epsilon\epsilon^T\right]$$

$$\text{Var}(\hat{\beta}) = \text{diag}\left[(X^TX)^{-1}\bar{E}\left[\epsilon\epsilon^T\right]\right] = \text{diag}\left((X^TX)^{-1}\right) \cdot \boxed{\text{Var}(\epsilon)} \rightarrow \sigma^2$$

$$\boxed{Var(\hat{\beta}) = \left(\boxed{\sigma^2}\right) diag\left((X^T X)^{-1}\right)}$$

$$\hat{\beta} \sim N\left(\beta, \; \sigma^2 diag\left((X^T X)^{-1}\right)\right)$$

$$s^2 = \frac{1}{n-p} \sum_{i=1}^{n} \left(y_i - \hat{y}_i\right)^2$$

$$\frac{\hat{\beta} - \beta}{\sqrt{\sigma^2 diag\left((X^T X)^{-1}\right)}} \sim t_{n-p} \longleftarrow$$

# of parameters in the linear regression model

$$P\left(\hat{\beta} - t_{n-p, 1-\frac{\alpha}{2}} \sqrt{s^2 diag\left((X^T X)^{-1}\right)} \leq \beta \leq \hat{\beta} - t_{n-p, \frac{\alpha}{2}} \sqrt{s^2 diag\left((X^T X)^{-1}\right)}\right) = 1 - \alpha$$

If we use $\alpha = 0.01$, we will obtain 99% CI