

Simple linear Reg. model

X and Y

$$Y = \beta_0 + \beta_1 X. \quad \boxed{\quad}$$

- If the two (random) variables are probabilistically related, then for a fixed value of X, there is uncertainty in the value of the Y.

$$Y = \beta_0 + \beta_1 X + \varepsilon \quad \text{where } \varepsilon \leftarrow \text{a random variable.}$$

↳ here is random error

T	X
4	1
9.5	10

→ There are parameters β_0, β_1 and σ^2 such that for any fixed value of the independent variable X the dependent variable is a random variable related to X through the model equation

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

The quantity ε in the model equation is the error - a random variable assumed to be symmetrically distributed with

$$E(\varepsilon) = 0 \text{ and } \text{Var}(\varepsilon) = \sigma_\varepsilon^2 = \sigma^2$$

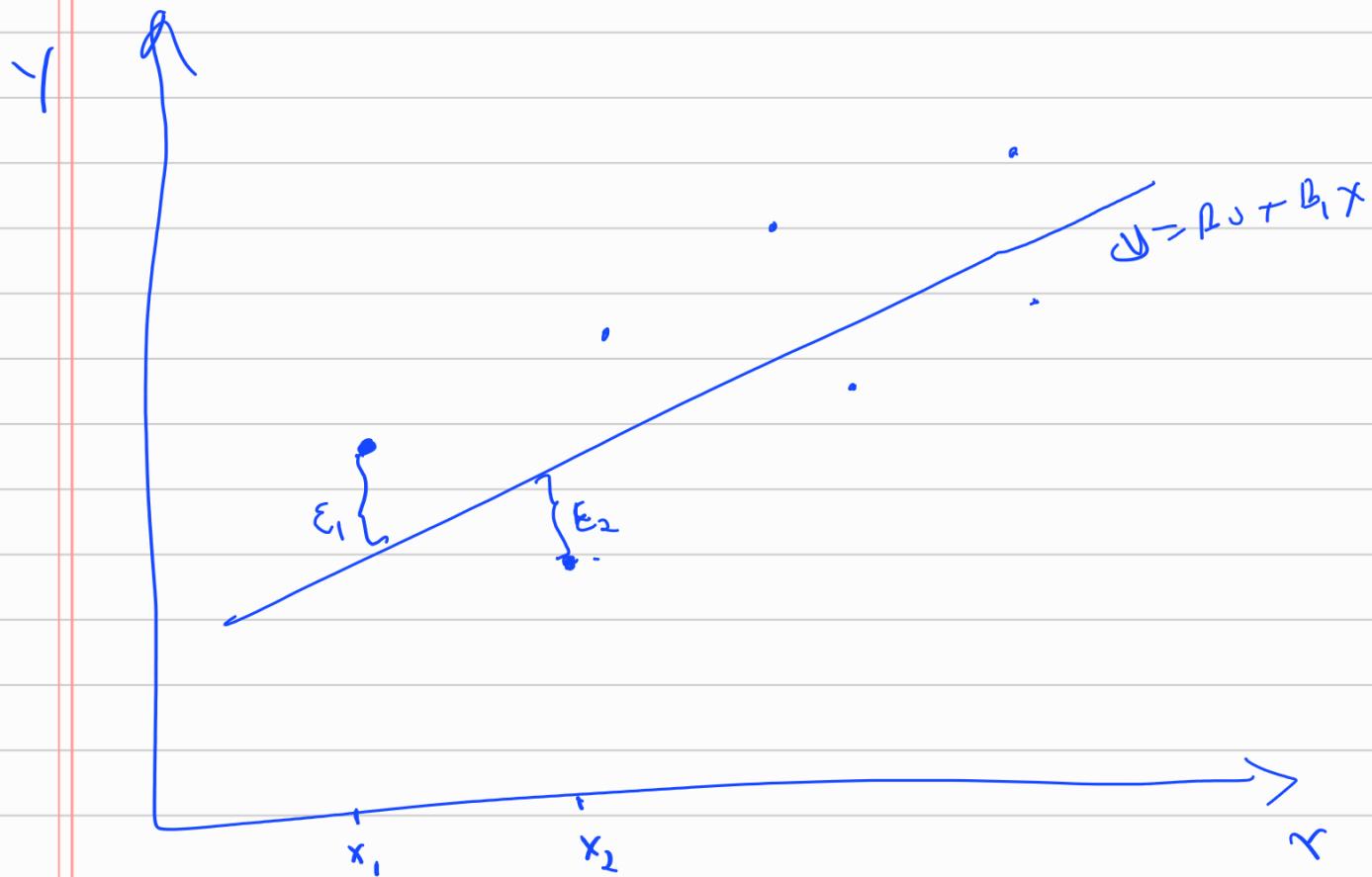
$$\varepsilon \sim N(0, \sigma^2) \approx \varepsilon \sim N(0, 1).$$

X : the independent, predictor or explanatory variable \rightarrow Not random.

Y : dependent / response variable. for fixed X , Y will be random variable.

ε : the random deviation or random error term for fixed x , ε will be random variable.

What exactly does ε do?



$\begin{matrix} \text{constant} \\ \text{coefficients} \end{matrix} \rightarrow \begin{matrix} \beta_0 \\ \beta_1 \end{matrix} \rightarrow \begin{matrix} \text{Vertical intercept of the model} \\ \text{Slope / Gradient} \end{matrix}$

Estimating the model parameters

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \text{ for } i=1, 2, \dots, n.$$

Least method: The sum of squared vertical distances from the point $(x_i, y_i), \dots, (x_n, y_n)$ to the line such that

$$\hat{f}(b_0, b_1) = \sum_{i=1}^n [y_i - (b_0 + b_1 x_i)]^2$$

- The first estimates of β_0 and β_1 , denoted by $\hat{\beta}_0$ and $\hat{\beta}_1$, are called the least squares estimates. - they give the values that minimize

$$f(b_0, b_1) =$$

$$\Rightarrow f'(b) = 0$$

$$\frac{\partial f(b_0, b_1)}{\partial b_0} = \sum 2(y_i - (b_0 + b_1 x_i))(-1) = 0$$

$$\frac{\partial f(b_0, b_1)}{\partial b_0} = \sum 2(y_i - (b_0 + b_1 x_i))(-1) = 0$$

$$\frac{\partial L(b_0, b_1)}{\partial b_1} = \sum_{i=1}^n 2(y_i - (b_0 + b_1 x_i))(-x_i) = 0$$

③

We need to solve

$$\sum_{i=1}^n 2(y_i - (b_0 + b_1 x_i))(-1) = 0 \quad \text{--- (i)}$$

$$\sum_{i=1}^n 2(y_i - (b_0 + b_1 x_i))(-x_i) = 0 \quad \text{--- (ii)}$$

$$\Rightarrow \sum_{i=1}^n (y_i - (b_0 + b_1 x_i)) = 0 \quad \text{--- (iii)}$$

$$\Rightarrow \sum_{i=1}^n x_i (y_i - (b_0 + b_1 x_i)) = 0 \quad \text{--- (iv)}$$

$$\left\{ \begin{array}{l} \sum_{i=1}^n y_i - n b_0 - b_1 \sum_{i=1}^n x_i = 0 \\ \sum_{i=1}^n x_i y_i - b_0 \sum_{i=1}^n x_i - b_1 \sum_{i=1}^n x_i^2 = 0 \end{array} \right.$$

$$b_1 = \hat{B}_1 = \frac{s_{xy}}{s_{xx}} \quad \checkmark$$

$$S_{xy} = \sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n}$$

$$S_{xx} = \sum x_i^2 - \frac{(\sum x_i)^2}{n}$$

$$b_0 = \hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x} \quad \checkmark$$

In some =

$$\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{w(x_i - \bar{x})}{\text{Var}(x)}$$

$$= \frac{S_{xy}}{S_{xx}}$$

$$y = b_0 + b_1 x$$

b_0 = remain.

$$(w(x_i - \bar{x}) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n})$$

Regression assumptions

- a) Linearity of the data: The relationship between the predictor (x) and the outcome (y) is assumed to be linear.
- b) Normality of residuals: The residual errors are assumed to be normally distributed $\text{Er} \sim N(0, \sigma^2)$
- c) Homogeneity of residual variance: The residuals are assumed to have a constant variance (homoscedasticity)
- d) Independence of the residual errors

Hypothesis test for linear models

- Hypothesis? A claim about pop parameters - Test the claim to establish whether it is true or not.

Types

- { Null hypothesis (H_0) ✓ maintains the status quo (neutral)
Alternative, (H_A, H_t, H_i) challenges

- Which one do we test?

$$H_0: \quad | \quad | \quad | \quad | \quad | \quad |$$

Decision: We reject H_0 (false) / Fail to reject
(true)

$$\Rightarrow \begin{array}{l} H_0: \mu = \mu_0 \\ H_A: \mu \neq \mu_0 \end{array} \quad \left. \begin{array}{l} \text{vs} \\ \text{Fail to reject} \end{array} \right\} \text{Two-tailed} \quad \left. \begin{array}{l} \text{Fail to reject} \\ \text{Do not reject} \end{array} \right\}$$

$$\Rightarrow \begin{array}{l} H_0: \mu = \mu_0 \\ H_A: \mu > \mu_0 \end{array} \quad \left. \begin{array}{l} \text{Fail to reject} \\ \text{Do not reject} \end{array} \right\} \text{one tail to right}$$

\hookrightarrow We fail to reject / We do not reject H_0 .

"We commit errors": How many errors do we commit while conducting testing hypothesis

\Rightarrow a) Type I error - Reject H_0 (H_0 true)

5) Type II error - Fail to | H₀ false

$$\lambda = P(\text{Type I error})$$

$$= P(\text{Reject } H_0 \mid H_0 \text{ true})$$

→ The significance level or size of the test

$$\alpha = 0.05, 0.02, 0.01,$$

$$\Rightarrow \text{Power of the test} \Rightarrow \beta = 0.25$$

$$\text{Power of the test} = 1 - \beta = \underline{99.75}$$

$$0.25 \rightarrow \\ \cdot \beta$$

Linear model

$$Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \hat{\beta}_2 X_{i2} + \dots + \hat{\beta}_{(p-1)} X_{i(p-1)} + \epsilon_i$$

$$\text{where } \epsilon_i \sim N(0, \sigma^2)$$

Test on a single parameter

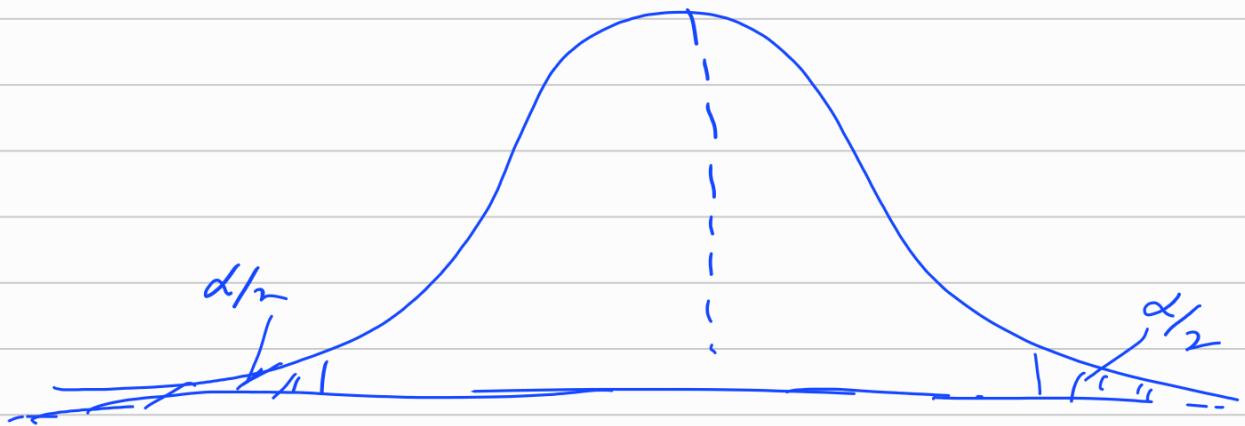
$$H_0: \beta_j = b \quad \text{vs} \quad H_1: \beta_j \neq b \quad \checkmark$$

- Significance level 100α%, b → constant
 $\alpha = 0.05$ or 0.01 or 0.02.

$$T = \frac{\hat{\beta}_j - b}{\text{SE}(\hat{\beta}_j)} \geq t_{(n-p), \alpha/2}$$

$$\therefore S \in (\hat{\beta}_j) = \sqrt{\text{Var}(\hat{\beta}_j)} \rightarrow \text{Standard error}$$

$t_{(n-p), \frac{\alpha}{2}}$ → Two-tailed test.
 Degrees of freedom



Test Statistics

$$H_0: \beta_j = 0 \quad \text{vs} \quad H_1: \beta_j \neq 0$$

$$|T| = \left| \frac{\hat{\beta}_j}{S E(\hat{\beta}_j)} \right| \sim t_{(n-p)}$$

a $\text{100}(1-\alpha)\%$ confidence interval for β_j

$$\left(\hat{\beta}_j - t_{(n-p), \frac{\alpha}{2}} S E(\hat{\beta}_j), \hat{\beta}_j + t_{(n-p), \frac{\alpha}{2}} S E(\hat{\beta}_j) \right)$$

constraint = β_0 .

- H₀: $\beta_0 = 0$ vs $\beta_0 \neq 0$.
- $\alpha = 0.05 \cdot n$

a) P-value:

Reject H₀: when p-value is less than α
· (level of significance).

$$P\text{-value} = 0.000$$

Decision:

Conclusion:

- Confidence interval

Test for the existence of regression model