

Simple Linear Regression (Review)

Regression analysis is a statistical technique used to describe relationships among variables

the simplest case : one explanatory variable observations:

SLR Model : $y_i = \beta_0 + \beta_1 x_i + \epsilon_i, i = 1, \dots, n$ $(x_i, y_i), i = 1, \dots, n$

$y_i \in \mathbb{R}$ is the real-value response for the i -th observation (dependent variable)

β_0 and β_1 (parameters) - regression coefficients

$x_i \in \mathbb{R}$ is the predictor or explanatory variable (independent variable)

ϵ_i random variables (error term) ~~unobserved~~ unobserved

Assumption (random error) : $E(\epsilon_i) = 0, \text{Var}(\epsilon_i) = \sigma^2, \forall i$

$\text{Cov}(\epsilon_i, \epsilon_j) = 0 \forall i \neq j$ (non-correlated)

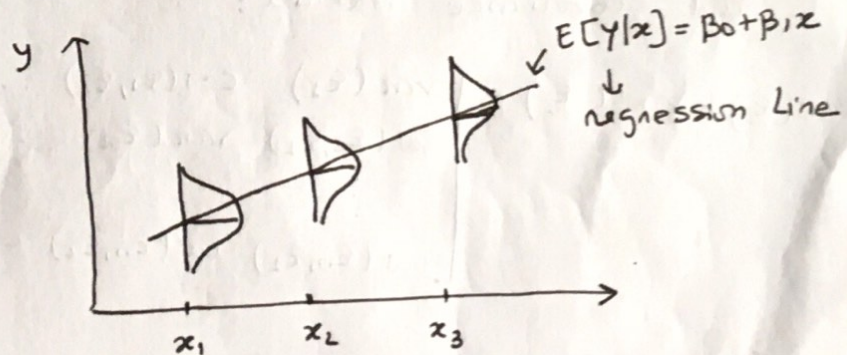
usually we also assume that $\epsilon_i \sim N(0, \sigma^2)$ Gaussian error term iid (independent and identically distributed)

- the model is simple because we have only one predictor
- the model is linear because y_i is a linear function of the parameters β_0 and β_1
- the model is a regression model because we are modeling a response variable (y) as a function of a predictor variable (x)

$E[y|x_i] = E[\beta_0 + \beta_1 x_i + \epsilon_i] = \beta_0 + \beta_1 x_i$

$\text{Var}(y|x_i) = \text{Var}(\beta_0 + \beta_1 x_i + \epsilon_i) = \text{Var}(\epsilon_i) = \sigma^2, \forall i$ homogeneity of variance

$\Rightarrow (y|x_i) \sim N(\beta_0 + \beta_1 x_i; \sigma^2)$ indep



$$E[Y|x] = \boxed{\beta_0 + \beta_1 x} \text{ regression Line}$$

β_0 - Intercept
 β_1 - Slope
 } unknown parameters
and $\text{var}(Y|x) = \text{var}(\epsilon) = \sigma^2$?
↑
unknown

Matrix notation:

$$\begin{aligned}
 Y_1 &= \beta_0 + \beta_1 x_1 + \epsilon_1 \\
 Y_2 &= \beta_0 + \beta_1 x_2 + \epsilon_2 \\
 &\vdots \\
 Y_n &= \beta_0 + \beta_1 x_n + \epsilon_n
 \end{aligned}$$

$$\boxed{\underset{\sim}{Y} = \underset{\sim}{X} \underset{\sim}{\beta} + \underset{\sim}{\epsilon}}$$

Notation: Matrix and vectors in bold face letters

$$\underset{\sim}{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \text{ (n x 1) vector of response random variables}$$

$$\underset{\sim}{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \text{ (n x 2) matrix called the design matrix}$$

$$\underset{\sim}{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \text{ (2 x 1) vector of unknown parameters}$$

$$\underset{\sim}{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix} \text{ (n x 1) vector of random errors}$$

$$\underset{\sim}{I} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & 1 \end{bmatrix} \rightarrow \text{Identity matrix}$$

$\underset{\sim}{\epsilon} \sim N_n(\underset{\sim}{0}, \sigma^2 \underset{\sim}{I})$ is an n-dimensional normal distribution with:

$$E(\underset{\sim}{\epsilon}) = E \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix} = \begin{bmatrix} E(\epsilon_1) \\ E(\epsilon_2) \\ \vdots \\ E(\epsilon_n) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \underset{\sim}{0}_{n \times 1}$$

and Covariance matrix:

$$\text{cov}(\underset{\sim}{\epsilon}) = \text{var}(\underset{\sim}{\epsilon}) = \begin{bmatrix} \text{var}(\epsilon_1) & \text{cov}(\epsilon_1, \epsilon_2) & \dots & \text{cov}(\epsilon_1, \epsilon_n) \\ \text{cov}(\epsilon_2, \epsilon_1) & \text{var}(\epsilon_2) & \dots & \text{cov}(\epsilon_2, \epsilon_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(\epsilon_n, \epsilon_1) & \text{cov}(\epsilon_n, \epsilon_2) & \dots & \text{var}(\epsilon_n) \end{bmatrix} = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix} \\
 = \sigma^2 \underset{\sim}{I}_{n \times n}$$

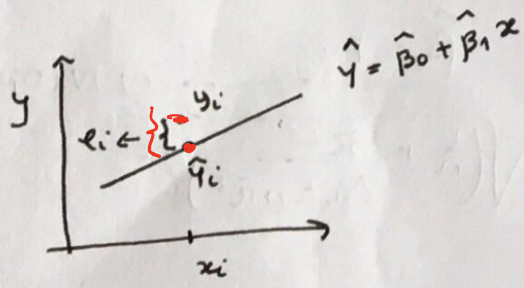
$$\begin{aligned}
 \text{cov}(\epsilon_i, \epsilon_j) &= E(\epsilon_i \epsilon_j) - E(\epsilon_i) E(\epsilon_j) \\
 \text{cov}(\epsilon_i, \epsilon_i) &= E(\epsilon_i \epsilon_i) - E(\epsilon_i) E(\epsilon_i) \\
 &= E(\epsilon_i^2) - E^2(\epsilon_i) = \text{var}(\epsilon_i)
 \end{aligned}$$

* Estimation of SLR model:

$$E[Y|x] = \beta_0 + \beta_1 x$$

$$\hat{E}[Y|x_i] = \hat{y}_i = \hat{\mu}_{Y|x_i} = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

$\min x_i \leq x_i \leq \max x_i$



Residuals: $y_i - \hat{y}_i = e_i = \varepsilon_i$

Least square Method: find $\hat{\beta}_0$ and $\hat{\beta}_1$ that minimize the sum of square residuals (it is also called the sum of squared errors, SSE)

$$\min_{\hat{\beta}_0, \hat{\beta}_1 \in \mathbb{R}} SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

$$\begin{cases} \frac{\partial SSE}{\partial \hat{\beta}_0} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \\ \frac{\partial SSE}{\partial \hat{\beta}_1} = -2 \sum_{i=1}^n x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \end{cases} \quad \text{Normal equations}$$

Solution:

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}$$

the least square estimate of σ^2 is $\hat{\sigma}^2 = \frac{SSE}{n-2} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2} = \dots =$

$$= \frac{1}{n-2} \left[\left(\sum_{i=1}^n y_i^2 - n \bar{y}^2 \right) - \hat{\beta}_1^2 \left(\sum_{i=1}^n x_i^2 - n \bar{x}^2 \right) \right]$$

↓
computational formula

* Inference in SLR

$\varepsilon_i \sim N(0, \sigma^2)$ i.i.d. $\Rightarrow y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$, y_i indep

Parameter	Estimator	Expected value:	Variance
Intercept: β_0	$\hat{\beta}_0$	$E(\hat{\beta}_0) = \beta_0$	$\text{Var}(\hat{\beta}_0) = \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum x_i^2 - n \bar{x}^2} \right) \sigma^2$
Slope: β_1	$\hat{\beta}_1$	$E(\hat{\beta}_1) = \beta_1$	$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum x_i^2 - n \bar{x}^2}$
$E[Y x_0] = \beta_0 + \beta_1 x_0$	$\hat{\mu}_{Y x_0} = \hat{\beta}_0 + \hat{\beta}_1 x_0$	$E[\hat{\mu}_{Y x_0}] = \beta_0 + \beta_1 x_0$	$\text{Var}(\hat{\mu}_{Y x_0}) = \left(\frac{1}{n} + \frac{(\bar{x} - x_0)^2}{\sum x_i^2 - n \bar{x}^2} \right) \sigma^2$

* Inference on β_0 :

Pivotal quantity: $\hat{\beta}_0 \sim N(\beta_0; (\frac{1}{n} + \frac{\bar{x}^2}{\sum x_i^2 - n\bar{x}^2}) \sigma^2)$

$$\frac{\hat{\beta}_0 - \beta_0}{\sqrt{(\frac{1}{n} + \frac{\bar{x}^2}{\sum x_i^2 - n\bar{x}^2}) \sigma^2}} \sim N(0, 1)$$

unknown σ^2

$$T = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{(\frac{1}{n} + \frac{\bar{x}^2}{\sum x_i^2 - n\bar{x}^2}) \hat{\sigma}^2}} \sim t_{(n-2)}$$

estimator

* β_1 : Pivotal $T = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\hat{\sigma}^2}{\sum x_i^2 - n\bar{x}^2}}} \sim t_{(n-2)}$

* $\mu_{Y|x_0}$: Pivotal variable: $T = \frac{(\hat{\beta}_0 + \hat{\beta}_1 x_0) - (\beta_0 + \beta_1 x_0)}{\sqrt{(\frac{1}{n} + \frac{(\bar{x} - x_0)^2}{\sum x_i^2 - n\bar{x}^2}) \hat{\sigma}^2}} \sim t_{(n-2)}$
 quantity:
 $= E[Y|x_0] = \beta_0 + \beta_1 x_0$

* Prediction Interval for a new or future observation of Y (Y_0)

$$(Y|x_0) \sim N(\beta_0 + \beta_1 x_0; \sigma^2)$$

$$Y|x_0 = Y_0 = \beta_0 + \beta_1 x_0 + \epsilon = E[Y|x_0] + \epsilon$$

$$\hat{Y}|x_0 = \hat{E}[Y|x_0] = \hat{\beta}_0 + \hat{\beta}_1 x_0$$

$$\hat{Y}_0 - Y_0 = (\hat{\beta}_0 + \hat{\beta}_1 x_0 - \beta_0 - \beta_1 x_0 - \epsilon)$$

$$E(\hat{Y}_0 - Y_0) = 0$$

$$\text{var}(\hat{Y}_0 - Y_0) = \text{var}(\hat{Y}_0) + \text{var}(Y_0) = \text{var}(E[Y|x_0]) + \sigma^2$$

$$= \sigma^2 \left[\frac{1}{n} + \frac{(\bar{x} - x_0)^2}{\sum x_i^2 - n\bar{x}^2} + 1 \right]$$

Pivotal quantity: $T = \frac{(\hat{Y}_0 - Y_0) - E(\hat{Y}_0 - Y_0)}{\sqrt{\text{var}(\hat{Y}_0 - Y_0)}} =$

$$= \frac{\hat{Y}_0 - Y_0}{\sqrt{\hat{\sigma}^2 \left(1 + \frac{1}{n} + \frac{(\bar{x} - x_0)^2}{\sum x_i^2 - n\bar{x}^2} \right)}} \sim t_{(n-2)}$$



Regression Sum-of-Squares

- sum of squares total: $SST = \sum_{i=1}^n (y_i - \bar{y})^2$
- sum of squares Regression: $SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$
- sum of squares Error (Squares of residuals): $SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$

Campus do Taguspark

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\underbrace{y_i - \hat{y}_i}_a + \underbrace{\hat{y}_i - \bar{y}}_b)^2$$

$$= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y})$$

= 0 under the validity of the model

$SST = SSR + SSE$

Coefficient of determination

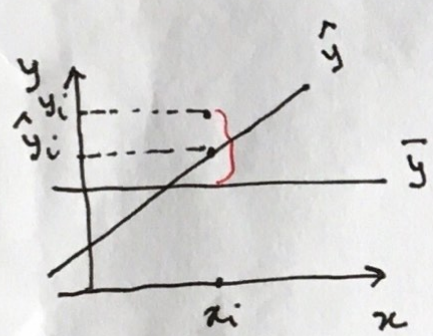
$$R^2 = \frac{SSR}{SST} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

gives the amount of variation in y 's that is explained by the linear relationship with x

$0 \leq R^2 \leq 1$

$R^2 \rightarrow 1$ (model good for this data)

$R^2 \rightarrow 0$ (not good model for this data)



$$(y_i - \bar{y}) = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$$

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\underbrace{y_i - \hat{y}_i}_a + \underbrace{\hat{y}_i - \bar{y}}_b)^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y})$$

= 0 under the validity of the model

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$SST = SSE + SSR$

(total) (residuals) (regression)

Example: Pizza Data

the owner of a Pizza restaurant ^{chefe} believes that if the restaurant is located near a college campus then there is a linear relationship between sales and the size of the student population.

Sample of 10 Pizza restaurants located near college campus

Population (1000s): x	2	6	8	8	12	16	20	20	22	26
Sales (\$1000s): y	58	105	88	118	117	137	157	169	149	202

$\sum_{i=1}^{10} x_i = 140$; $\sum_{i=1}^{10} y_i = 1300$; $\sum_{i=1}^{10} x_i^2 = 2528$; $\sum_{i=1}^{10} y_i^2 = 184730$; $\sum_{i=1}^{10} x_i y_i = 21040$

a) Estimate the least-squares regression line

$$\bar{x} = 14 ; \bar{y} = 130 ; \sum_{i=1}^{10} x_i^2 - n\bar{x}^2 = 2528 - 10 \times 14^2 = 568$$

$$\sum_{i=1}^{10} x_i y_i - n\bar{x}\bar{y} = 21040 - 10 \times 14 \times 130 = 2840$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{10} x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^{10} x_i^2 - n\bar{x}^2} = \frac{2840}{568} = 5$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 130 - 5 \times 14 = 60 \quad \hat{y} = 60 + 5x$$

b) $\hat{E}[Y|x=9] = 60 + 5 \times 9 = 105$; $\hat{E}[Y|x=27] \rightarrow$ extrapolation

c) Estimate the variance of the sales variable

$$\hat{\text{var}}(Y) = \hat{\sigma}^2 = \frac{1}{8} \left[\left(\sum_{i=1}^{10} y_i^2 - n\bar{y}^2 \right) - 5^2 \times 568 \right] = \frac{15730 - 5^2 \times 568}{8} = \frac{1530}{8}$$

$$= 191.25$$

$$\sum_{i=1}^{10} y_i^2 - n\bar{y}^2 = 184730 - 10 \times 130^2 = 15730$$

d) Estimate the residual of the second observation

$$y_2 = 105 \quad \hat{y}_2 = 60 + 5 \times 6 = 90 \quad e_2 = 105 - 90 = 15$$

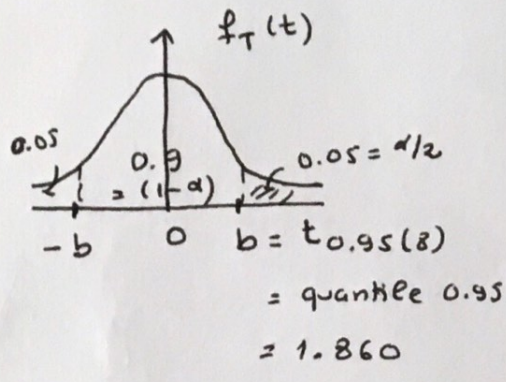
e) construct a 90% CI for β_1

$$CI_{90\%}(\beta_1) = ?$$

Pivotal variable: $T = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^{10} x_i^2 - n\bar{x}^2}}} \sim t(8)$



$P(a \leq T \leq b) = 0.9$
 symmetrical interval : $a = -b$



$P(-1.860 \leq \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\hat{\sigma}^2}{\sum x_i^2 - n\bar{u}^2}}} \leq 1.860) = 0.95$

$P(\hat{\beta}_1 - 1.860 \times \sqrt{\frac{\hat{\sigma}^2}{\sum x_i^2 - n\bar{u}^2}} \leq \beta_1 \leq \hat{\beta}_1 + 1.860 \times \sqrt{\frac{\hat{\sigma}^2}{\sum x_i^2 - n\bar{u}^2}}) = 0.9$
 Random variable

CI(β_1) = [3.921; 6.079]

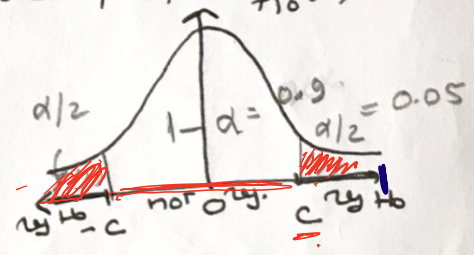
90% null hypothesis Alternative hypothesis

d) ~~test~~ test: $H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$ significance level 10%.
 (not significant) (significant)
 very important test: If the null hypothesis is not reject the model can be use since the variable x doesn't explain the expected value of y .
 test the significance of the regression model
 two sided test

Pivotal variable:
 quantity: $T = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\hat{\sigma}^2}{\sum x_i^2 - n\bar{u}^2}}} \sim t(8)$

test statistic: under the validity of H_0 $f_{T_0}(t)$

$t_0 = T|_{H_0} = \frac{\hat{\beta}_1 - 0}{\sqrt{\frac{\hat{\sigma}^2}{\sum x_i^2 - n\bar{u}^2}}} \sim t(8)$



α = significance level
 $= P(\text{rej } H_0 | H_0 \text{ is true}) = P(|T_0| > c) = 0.1$

$\Rightarrow c = t_{0.95}(8) = F(0.95) = 1.860$

critical region: CR = $]-\infty, -1.860[\cup]1.860; +\infty[$

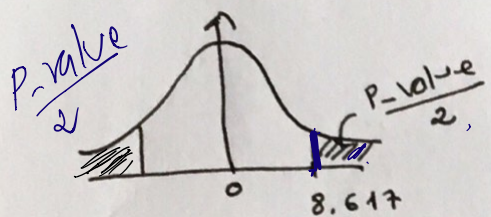
test statistic observed value: $t_0 = \frac{5 - 0}{\sqrt{\frac{191.25}{568}}} \approx 8.617$

Decision: Since $t_0 \in CR$, for $\alpha = 0.1$, reject H_0 , i.e. ~~with~~ this sample it looks like that the size of the student population explain ~~linearly~~ the expected value of sale (there is a linear relationship between sale, and the student population size)

Alternative: not fix significance level

P-value of the test:

Statistic observed value: $t_0 \approx 8.617$



$$P\text{-value} = 2(P(T_0 \geq 8.617)) = 2(1 - F_{T_0}(8.617)) =$$

table: $F_{t(8)}(5.041) = 0.9995$

$$0.9995 < F_{t_8}(8.617) < 1 \Rightarrow 0 < 1 - F_{T_0}(8.617) < (1 - 0.9995)$$

$$\Rightarrow 0 < P\text{-value} < 0.001$$

$\Rightarrow \forall \alpha \geq 0.001$ rej $H_0 \Rightarrow$ rej H_0 to use

e) test: $H_0: \beta_0 = 75$ vs $H_1: \beta_0 < 75$ (one sided test)

use $\alpha = 0.05$

Left tailed test

test statistic: $T_0 = \frac{\hat{\beta}_0 - 75}{\sqrt{\widehat{\text{var}}(\hat{\beta}_0)}} \sim t(8)$

$$\widehat{\text{var}}(\hat{\beta}_0) = 85.1197$$

observed value $t_0 = \frac{60 - 75}{\sqrt{85.1197}} \approx -1.626$

critical region:

$$C.R =]-\infty, -1.860[$$

since $t_0 \notin C.R.$ not rej. H_0

