

Linear Regression

Contingency Table Tests allow us to explore association between two categorical variables:

$$R_y = \{\text{Blue, Green, Red, male, Turkish}\}$$

Regression analysis also allow us to explore association between two variables but instead two numeric variables:

$$Y = \{2.18, 4.28, 7.92, \dots\}$$

$$X = \{1.1, 1.2, 1.4, \dots\}$$

R.V.

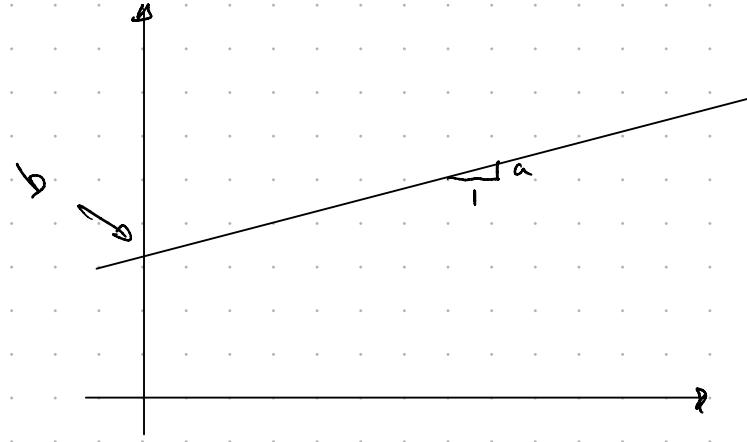
Discrete: Contingency tables

Continuous: Regression

$$y = ax + b$$

Dependent

Independent



$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots$$

The association is called correlation between

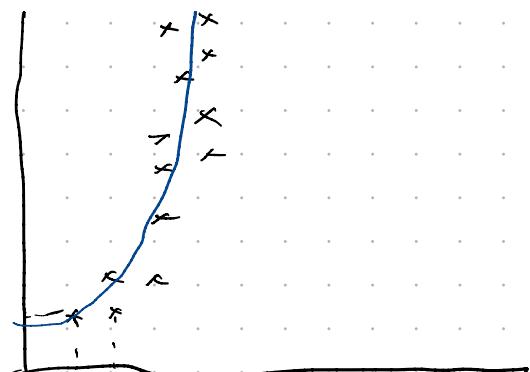
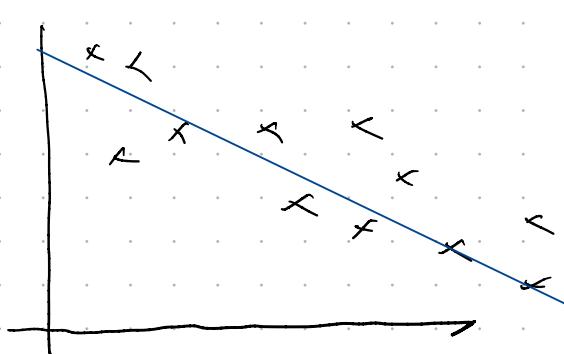
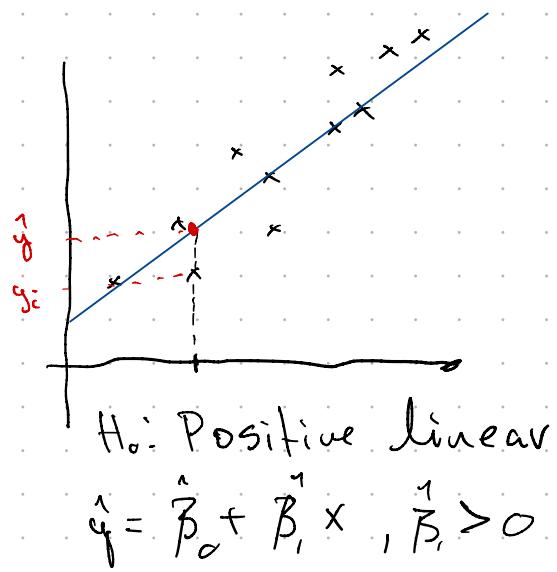
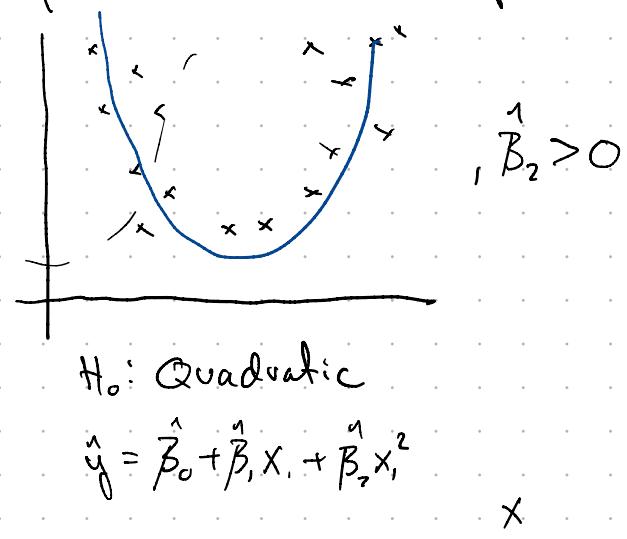
X and y :

$$X = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 & x_5 \\ p_1 & 16 & 47 & 122 & \cdot \\ p_2 & 18 & 72 & \cdot & \cdot \\ p_3 & 20 & 94 & \cdot & \cdot \\ p_4 & 15 & 31 & \cdot & \cdot \end{bmatrix}$$

$$X = \begin{bmatrix} 16 & 87 & 122 \\ 18 & 72 & \cdot \\ 20 & 94 & \cdot \\ 15 & 31 & \cdot \end{bmatrix}$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{bmatrix}, X \cdot \beta$$

① Scatter Plot of X and Y to visually inspect relationship.



2. Remove outliers. In regression this is done qualitatively

3. Determine regression equations, i.e. estimate β_0 and β_1 :

$$y = \beta_0 + \beta_1 x + \varepsilon$$

$$\left. \begin{array}{l} \hat{\beta}_0 = ? \\ \hat{\beta}_1 = ? \end{array} \right\} \quad \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$EY = \beta_0 + \beta_1 \cdot E[X] + E(\varepsilon)$$

$$= \beta_0 + \beta_1 \cdot E[X]$$

$$\beta_0 = EY - \beta_1 \cdot E[X]$$

$$\text{Cov}(X, Y) = (X, \beta_0 + \beta_1 X + \varepsilon)$$

$$= \beta_0 \text{Cov}(X, 1) + \beta_1 \text{Cov}(X, X) + \text{Cov}(X, \varepsilon)$$

$$= 0 + \beta_1 \cdot \text{Var}(X) + 0$$

$$= \beta_1 \cdot \text{Var}(X)$$

$$\beta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}, \quad \beta_0 = EY - \beta_1 \cdot E[X]$$

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}, \quad \text{Cov}(X, Y) = \frac{1}{n-1} \cdot \sum (x_i - \bar{x}) \cdot (y_i - \bar{y})$$

$$\bar{y} = \frac{y_1 + y_2 + y_3 + \dots + y_n}{n}, \quad \text{Var}(X) = \frac{1}{n-1} \cdot \sum (x_i - \bar{x})^2$$

$$\Rightarrow \text{num} = S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{den} = S_{xx} = \sum (x_i - \bar{x})^2$$

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

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

4. Check assumptions that errors are normally distributed (normal probability plot)

Residual(error)

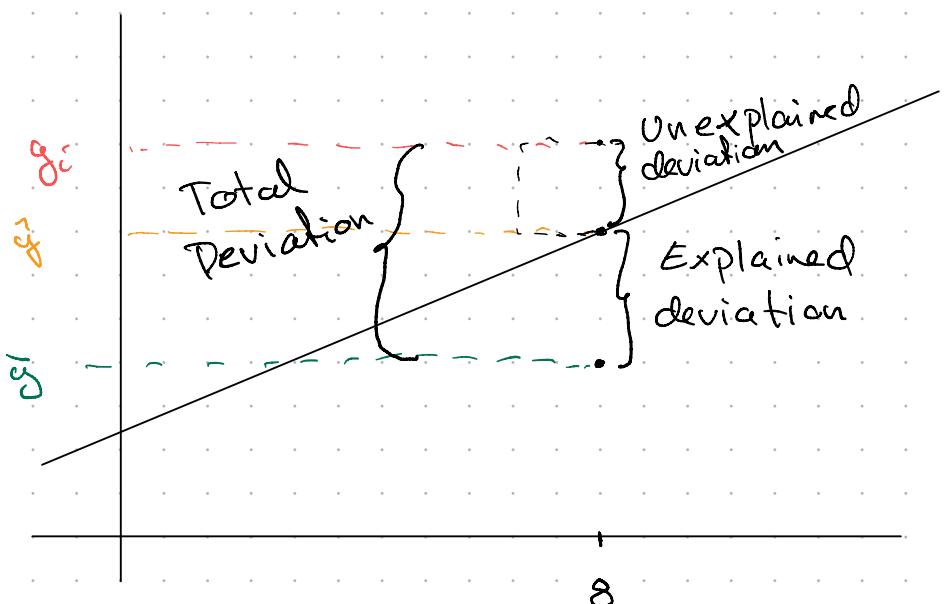
$$e_i = y_i - \hat{y}_i$$

$y_i = 9 \rightarrow$ Observed

$\hat{y}_i = 7 \rightarrow$ Predicted

$\bar{y} = 4 \rightarrow$ average

$\lambda = 8.$



Total: $(y_i - \bar{y})$

Unexp: $(y_i - \hat{y}_i)$

expl: $(\hat{y}_i - \bar{y})$

Total sum of Squares:

$$SS_T = \sum (y_i - \bar{y})^2$$

Regression sum of Squares:

$$SS_R = \sum (\hat{y}_i - \bar{y})^2$$

Residual sum of squares:

$$SS_E = \sum (y_i - \hat{y}_i)^2 \leftarrow \text{Minimize this!}$$

S. Assess adequacy of model:

a) $H_0: \beta_1 = 0$

$H_1: \beta_1 \neq 0$

Test statistic: $T_0 = \frac{\hat{\beta}_1}{\sqrt{\hat{\sigma}^2 / S_{xx}}}, \hat{\sigma}^2 = \frac{SSE}{n-2}$

b) Determine correlation:

1) $r = \frac{\sum (Z_x \cdot Z_y)}{n-1}$ z scores of all x and y

2)

⋮

4) $r = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sqrt{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2}} = \frac{S_{xy}}{\sqrt{S_{xx} \cdot S_{yy}}}$

$|r| > 0.6$: good

$|r| \leq 0.8$: High

$0.4 \leq |r| \leq 0.6$: OK

c) Find Correlation of Determination:

$$r^2 = \text{square of } r \quad \left. \begin{array}{l} \text{amount of variance} \\ \text{my model is able} \end{array} \right\}$$

$$r^2 = \frac{SSE}{SST} \left. \begin{array}{l} \leftarrow \text{Expl} \\ \leftarrow \text{Total} \end{array} \right\}$$

to explain.

6) Find confidence intervals
for β_0 and β_1

