# Linear Regression

Contingency Table Tests allow to explore association between two categorical variable

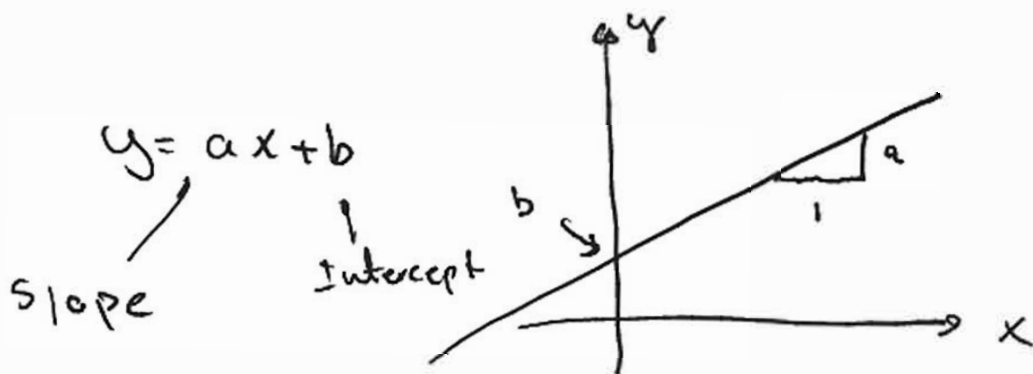$$\Omega_g = \{ Blue, Green, Red, male, Turkish \}$$

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

$$Y = \{ 2.18, 4.28, 7.92 \ldots \}$$
$$X = \{ 1.1, 1.2, 1.4, 1.8 \ldots \}$$

R.V.
  - Discrete → Contingency tables
  - continuous → Regression.

More specifically it is the type of the dependent variable that decides:



$y = ax + b$

Slope

Intercept

x is called the independent variable.
y is called the dependent or response variable.

This association is called the
correlation between X and Y.

Simple Linear Regression:

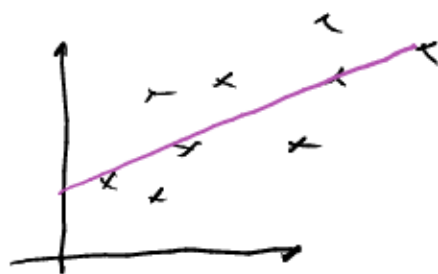One dependent and one independent.
→ Is there relation between x and y?
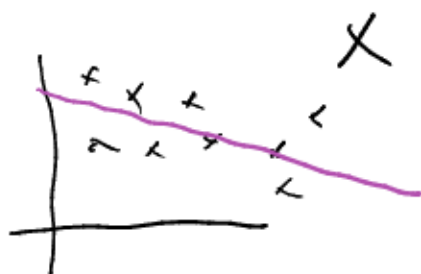 └ How strong is this relation?

A regression Analysis includes:

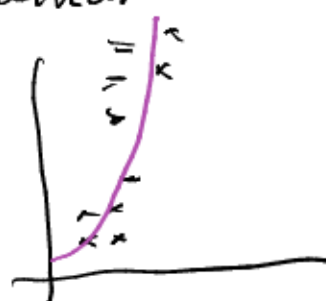1) Scatter plot of x and y to
visually inspect relationship



$H_o$: Quadratic

$H_o$: positive
linear

outlier

X

$H_o$: Negative
linear

$H_o$: exponential

2) Remove any outliers. In regression
this is qualitative assessment.
Visual inspection.

3) Determine the regression equation, i.e.
   estimate $\beta_0$ and $\beta_1$:

$$\boxed{y = \beta_0 + \beta_1 X + \varepsilon} \quad (y = b + ax)$$

$$\left.\begin{array}{l} \hat{\beta}_0 = ? \\[2mm] \hat{\beta}_1 = ? \end{array}\right\} \quad \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

$$Ey = \beta_0 + \beta_1 EX + E[\varepsilon]$$

$$= \beta_0 + \beta_1 EX$$

$$\beta_0 = EY - \beta_1 EX$$

Now we need $Cov(X,Y)$

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

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

$$= 0 + \beta_1 Var\,X + 0$$

$$\beta_1 = \frac{Cov(X,Y)}{Var\,X} \qquad , \quad \beta_0 = EY - \beta_1 EX$$

$$\bar{X} = \frac{X_1 + X_2 + \ldots X_n}{n}$$

$$\bar{y} = \frac{y_1 + y_2 + \ldots y_n}{n}$$

$$Cov(X,Y) = \frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$Var(x) = \frac{1}{n-1} \sum (x_i - \bar{x})^2$$

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

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

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

$\left.\rule{0pt}{40pt}\right\}$ Based on $Cov(x,y)$ / $Var(x)$

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

At exam I <u>will</u> ask for $S_{xy}$ and $S_{xx}$.

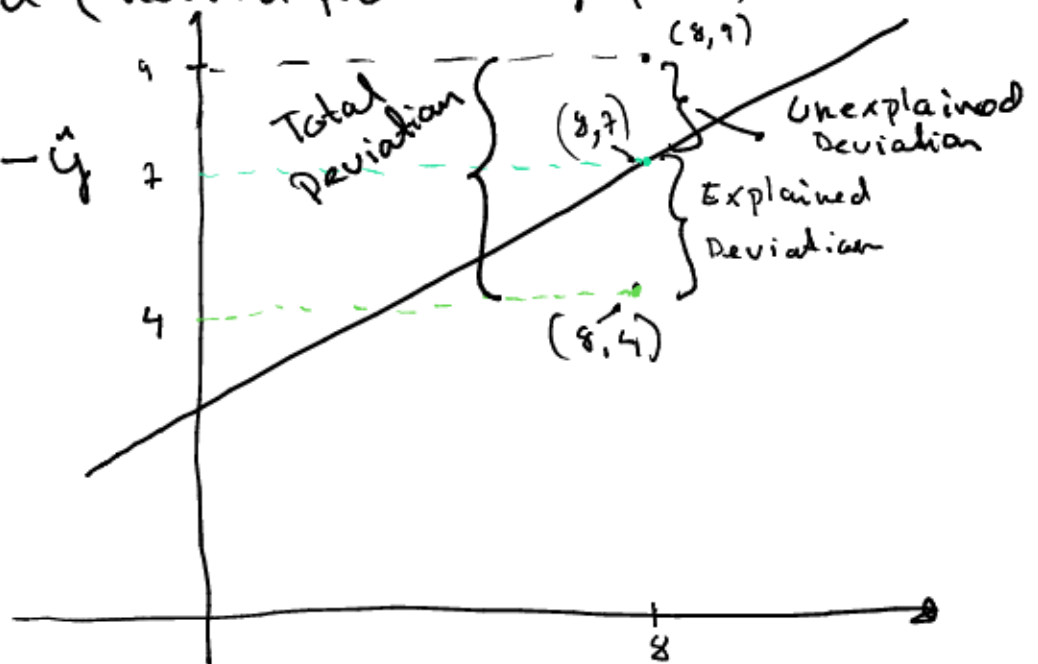4) Check assumption that errors are normally distributed (normal probability plot)

<u>Residual:</u>
$$e_i = y_i - \hat{y}$$

$y_i = 9 \rightarrow$ Observed

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

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



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

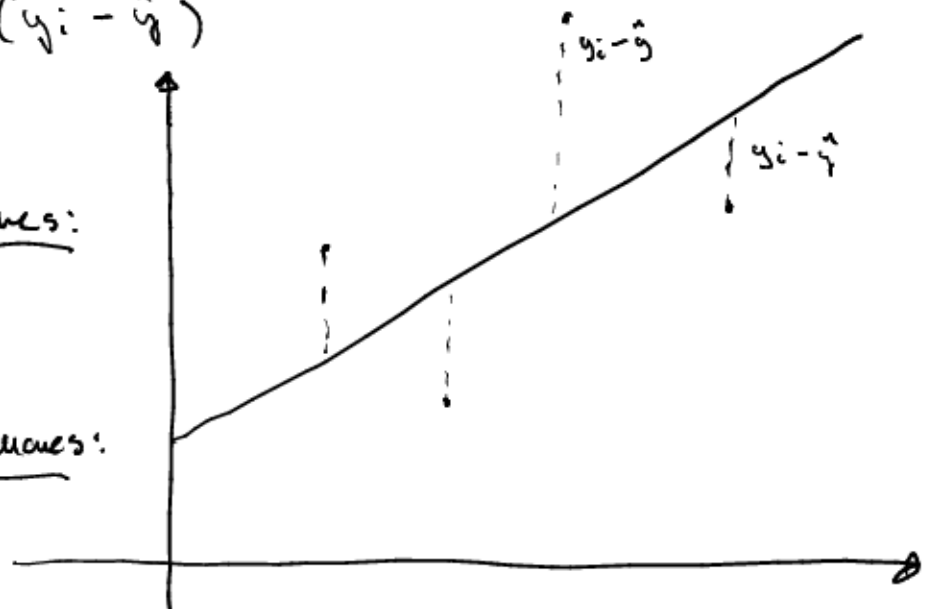Explained : $(\hat{y} - \bar{y})$

Unexplained : $(y_i - \hat{y})$

(Residual)

<u>Total Sum of squares:</u>

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

<u>Regression sum of squares:</u>

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

explained : $\hat{y} - \bar{y}$
Unexplained : $(y_i - \hat{y})$
(Residual)

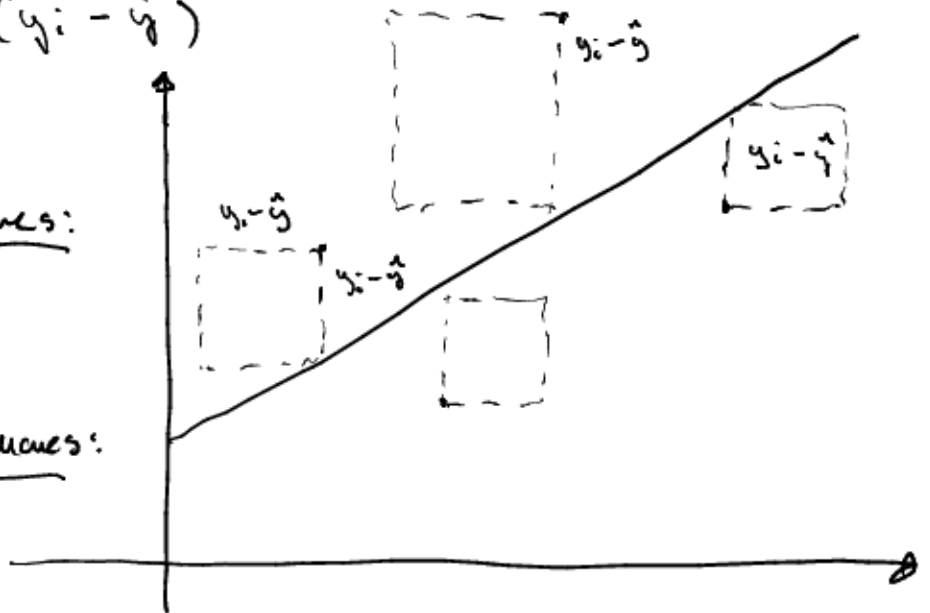Total Sum of squares:

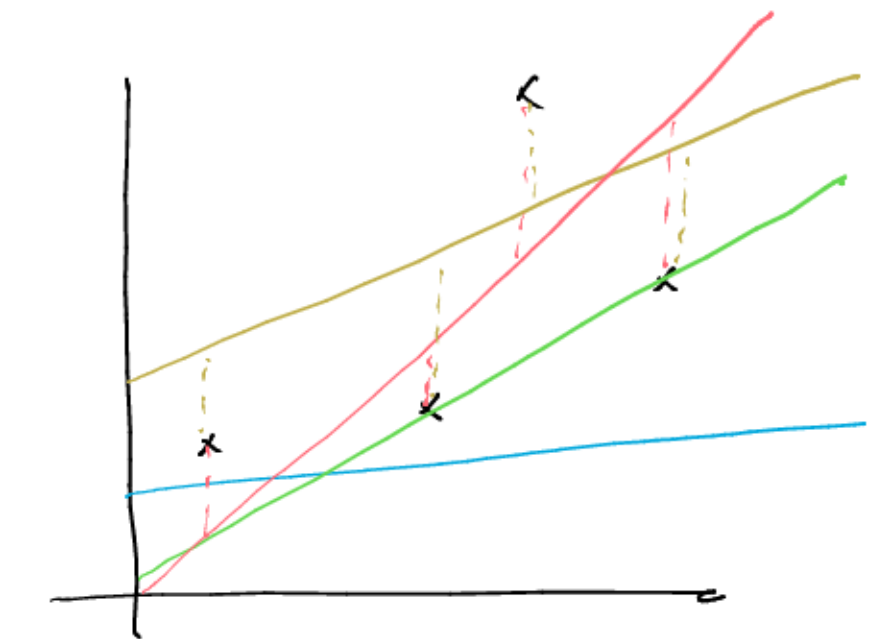$$SS_t = \sum (y_i - \bar{y})^2$$

Regression sum of squares:

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

Residual Sum of squares:

$$SS_E = \sum (y_i - \hat{y})^2$$

The goal of Linear Regression is to <u>minimize</u> the area of all of these squares

5) Assess adequacy of Model

    (a) Test: $H_0$: $\beta_1 = 0$

           $H_1$: $\beta_1 \neq 0$

       Test stat: $T_0 = \dfrac{\hat{\beta}_1}{\sqrt{\hat{\sigma}^2 / S_{xx}}}$ , $\hat{\sigma}^2 = \dfrac{SS_E}{n-2}$

    (b) Determine correlation

       1) $r = \dfrac{\Sigma (z_x \cdot z_y)}{n-1}$ , z scores for all sample values

       2) $r = \dfrac{n \cdot \Sigma xy - \Sigma x \cdot \Sigma y}{\sqrt{n \cdot \Sigma x^2 - (\Sigma x)^2} \sqrt{n \cdot \Sigma y^2 - (\Sigma y)^2}}$

       3) $r = \sqrt{\dfrac{SS_R}{SS_T}} = \sqrt{1 - \dfrac{SS_E}{SS_T}}$

       4) $r = \dfrac{\Sigma (y_i - \bar{y}) \cdot (x_i - \bar{x})}{\sqrt{\Sigma (x_i - \bar{x})^2 \cdot \Sigma (y_i - \bar{y})^2}} = \dfrac{S_{xy}}{\sqrt{S_{xx} \cdot SS_T}}$

    $r > |0.61| \rightarrow$ Good correlation

    $r > |0.81| \rightarrow$ High correlation

    (c) Find Correlation of Determination:

       $r^2 =$ square of $r$

       $r^2 = \dfrac{SS_R}{SS_T}$ ← Explained, ← Total $\Big\}$ $r^2$ tells me the amount of variance that my model is able to explain.

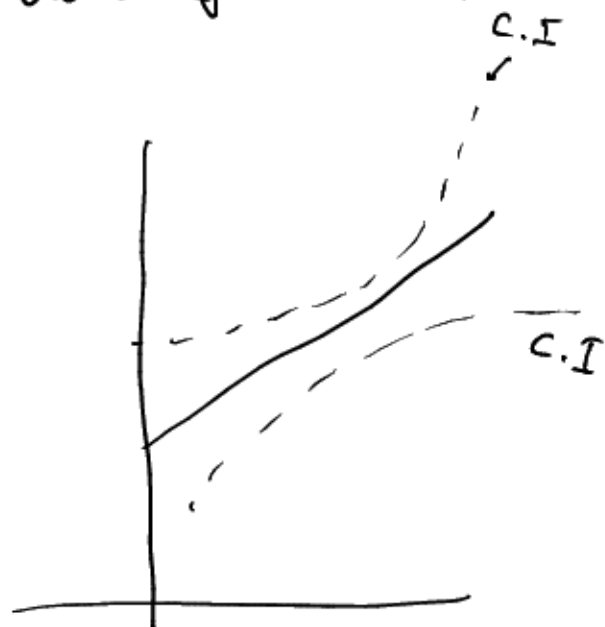       $r^2 = 0.78 \rightarrow r^2$ shows Prec

- $r^2$ tells me something about my model's predictive power (Proportion of explained error)
- $r$ tells me something about how well $x$ and $y$ co-vary

6) Find Confidence Intervals for slope and Intercept

Error for slope:

$$E_s = t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}$$

$$S(b_1) \rightarrow \text{standard error of slope}$$



C.I

C.I

Error for intercept:

$$E_I = t_{\alpha/2, n-2} \underbrace{\sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}}_{S(b_0) \rightarrow \text{standard error of intercept.}}$$

Prediction Intervals:

$$E_p = t_{\alpha/2, n-2} \cdot \sqrt{\hat{\sigma}^2 \left( 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)}$$

$x_0$ is the x-value we want to find $\hat{y}$ at.