

Simple Regression

Anis Rezgui
Mathematics Department
Carthage University - INSAT
Tunis - Tunisia

November 28, 2019

Table of contents

- 1 Introduction
- 2 Least square method
- 3 Interpretation
- 4 Inferential study

Introduction

The problem

Let x and y be two statistical variables:

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \text{and} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}.$$

- The correlation problem, in general, consists of looking for a possible relationship between x and y : $y = f(x) + \text{"residual"}$ that makes the residual part the **smallest** possible.
- The variable x is called "predictor" or the "independent" variable, y is called "response" or "dependent" variable.
- The "**linear**" correlation is when the relationship f is linear.

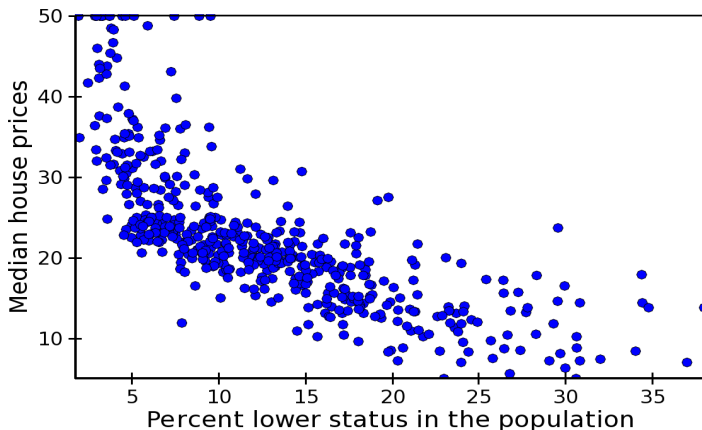
Introduction

The scattergram

The scattergram or "scatter-plot" is simply the set of points of coordinates

$$\{(x_i, y_i) : i = 1, \dots, n\}.$$

It gives a very useful a priori glance:

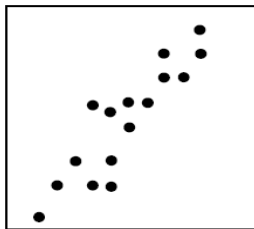


Introduction

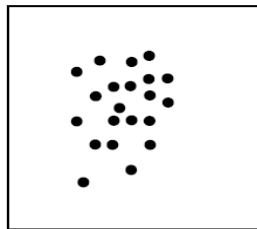
Different types of correlation



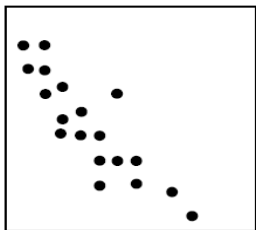
Strong positive correlation



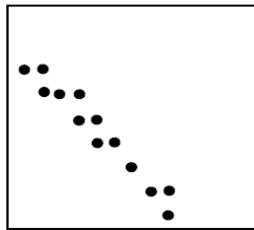
Moderate positive correlation



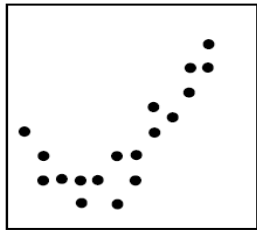
No correlation



Moderate negative correlation



Strong negative correlation



Curvilinear relationship

Least square method

The least square distance

Our problem can be reformulated as follows:

We look for β and α such that if $y_i^* = \beta x_i + \alpha$ the distance

$$\sum_{i=1}^n (y_i^* - y_i)^2 \text{ is minimal.} \quad (1)$$

- 1 The line $s = \beta t + \alpha$ is called "least square line" or "regression line".
- 2 β is called "slope" and α the "intercept" of the regression line.

Least square method

Vectorial formulation

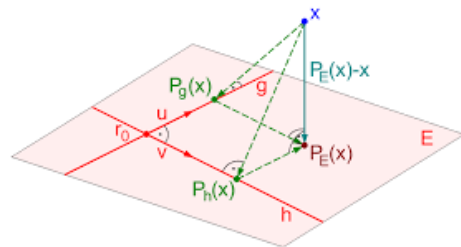
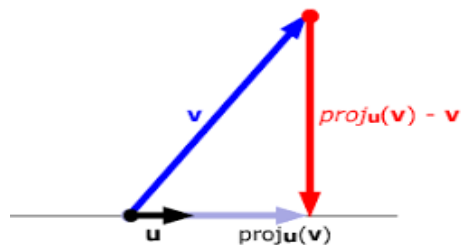
We look for a vector $y^* \in \mathbb{R}^n$ spanned by $\{x, \mathbf{1}\}$ i.e find out β and α such that $y^* = \beta x + \alpha \mathbf{1}$ minimizes the Euclidian norm $\|y - y^*\|_n^2$.

Theorem

$\|y - y^*\|_n^2$ is minimal if and only if

$y^* = \text{Orthogonal Projection of } y \text{ on the subspace } \text{span}\{x, \mathbf{1}\}.$

Least square method



Consequences

- $y - y^*$ is orthogonal to y^* .

$$y = y^* \oplus^\perp y - y^* \Rightarrow \|y\|_n^2 = \|y^*\|_n^2 + \|y - y^*\|_n^2 \quad (2)$$

- $\bar{y} = \overline{y^*}$ and (2)
 $\Rightarrow \|y - \bar{y}\|_n^2 = \|y^* - \bar{y}\|_n^2 + \|y - y^*\|_n^2$
- *Total sum of squares (tss) = fitted sum of squares (fss) + residual sum of squares (rss)*
- $\Rightarrow s_y^2 = s_{y^*}^2 + s_{y-y^*}^2$
Total variance = Explained variance + Residual variance

Least square method

Determination of the regression slope and intercept

We look for β and α so that $y - y^*$ is orthogonal to $\text{span}\{x, \mathbf{1}\}$ i.e

i. On one hand

$$\begin{aligned}y - y^* \perp x &\Rightarrow {}^t x(y - y^*) = 0 \\&\Rightarrow {}^t xy - \beta {}^t xx + \alpha n\bar{x} = 0 \\&\Rightarrow |x|^2 \beta - n\bar{x}\alpha = {}^t xy\end{aligned}\tag{3}$$

ii. on the other hand

$$\begin{aligned}y - y^* \perp \mathbf{1} &\Rightarrow {}^t \mathbf{1}(y - y^*) = 0 \\&\Rightarrow n\bar{y} - n\beta\bar{x} + n\alpha = 0 \\&\Rightarrow \alpha = \beta\bar{x} - \bar{y}\end{aligned}\tag{4}$$

Least square method

Finally

Combining (3) and (4)

$$\begin{cases} \beta = \frac{\overline{x * y} - \bar{x} * \bar{y}}{\overline{x^2} - \bar{x}^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ \alpha = \bar{y} - \beta \bar{x} \end{cases}$$

Linear Determination coefficient vs Correlation Coefficient

- 1 Set the Linear Determination Coefficient as

$$R^2_{y|x} = \frac{fss}{tss} = \frac{\sum_i (y_i^* - \bar{y})^2}{\sum_i (y_i - \bar{y})^2} \in [0, 1]$$

- 2 It represents the strongness of the linear correlation between y and x .
- 3 It can be read as the the rate of y explained linearly by x .

Linear Determination Coefficient vs Correlation Coefficient

Definition

The linear correlation coefficient of the statistical variables x and y is given by

$$r_{xy} = \frac{\text{Cov}(x, y)}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y})^2}} \quad (5)$$

$$= \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y})^2}} \quad (6)$$

$$= \frac{n}{n-1} \frac{\overline{x * y} - \bar{x} * \bar{y}}{s_x * s_y} \quad (7)$$

$$= \frac{\sum_i x_i y_i - \frac{(\sum_i x_i)(\sum_i y_i)}{n}}{\sqrt{\left(\sum_i x_i^2 - \frac{(\sum_i x_i)^2}{n}\right) \left(\sum_i y_i^2 - \frac{(\sum_i y_i)^2}{n}\right)}} \quad (8)$$

Very Important Remark

- i. The linear determination coefficient and the linear correlation coefficient are actually the same ! (Why?)

$$r_{xy}^2 = R_{y|x}^2. \quad (9)$$

- ii. The linear correlation coefficient can be seen as the **cosine** of the angle formed by the two vectors of \mathbb{R}^n , $x - \bar{x} * \mathbf{1}$ and $y - \bar{y} * \mathbf{1}$:

$$\begin{aligned} \cos(\widehat{x - \bar{x} * \mathbf{1}, y - \bar{y} * \mathbf{1}}) &= \frac{(x - \bar{x} * \mathbf{1})(y - \bar{y} * \mathbf{1})}{\|x - \bar{x} * \mathbf{1}\| \|y - \bar{y} * \mathbf{1}\|} \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \end{aligned}$$

Interpretation

Very Important Remarks

- ① Since the correlation coefficient is a cosine it satisfies:

$$|r_{xy}| \leq 1.$$

- ② If $|r_{xy}| = 1$ it means that we have a perfect linear correlation:

$$y_i = \beta x_i + \alpha, \text{ for all } i = 1, \dots, n.$$

- ③ If $r_{xy} = 0$ it means that there is no linear relationship between x and y .
- ④ In all cases we have

$$\beta = r_{xy} \times \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

Interpretation

So what ?

- ① If $|r_{xy}| \in [0, 25\%[$, there is no linear correlation !
- ② If $|r_{xy}| \in [25\%, 50\%[$, there is a moderate linear correlation.
- ③ If $|r_{xy}| \in [50\%, 75\%[$, there is a fair linear correlation.
- ④ If $|r_{xy}| \in [75\%, 100\%[$, there is good linear correlation.

Inferential study

Suppose that the two statistical series x and y came from two given random variables X and Y and let $\{X_i : i = 1 \cdots n\}$ and $\{Y_i : i = 1 \cdots n\}$ be two random samples of X and Y . Denote by

$$TSS = \sum_{i=1}^n (Y_i - \bar{Y})^2, \quad FSS = \sum_{i=1}^n (Y_i^* - \bar{Y})^2$$

and

$$RSS = \sum_{i=1}^n (Y_i - Y_i^*)^2$$

and consider the following statistic:

$$F = \frac{FSS}{RSS} (n - 2).$$

A theoretical result: The linear Gaussian model

Theorem

Suppose that X and Y are normally distributed and that they are independent. Then

$$\frac{FSS}{RSS}(n-2) = \mathcal{F}(1, n-2)$$

where $\mathcal{F}(1, n-2)$ is the Fisher distribution of degrees of freedom 1 and $n-2$.

Proposition and Definition

If $A \sim \chi_{d_1}^2$ and $B \sim \chi_{d_2}^2$ and are independent then

$$\frac{A/d_1}{B/d_2} \sim \mathcal{F}(d_1, d_2).$$

So what: How to use the latter Theorem ?

We consider the test: (H_0): $\beta = 0$ (which means that there is no linear relationship between X and Y) against (H_1): $\beta \neq 0$:

- we reject the hypothesis (H_0) and accept (H_1) with a confidence level of 95%, if:

$$F^* = \frac{f_{ss}}{r_{ss}}(n-2) \geq f_o$$

where f_o is such that:

$$\mathbb{P}\{\mathcal{F}(1, n-2) \leq f_o\} = 95\%$$

- or equivalently we reject H_0 and accept H_1 when

$$p - value = \mathbb{P}\left\{\mathcal{F}(1, n-2) \geq F^* = \frac{f_{ss}}{r_{ss}}(n-2)\right\} \ll 1$$

Confidence interval for the regression line

Let $T_{n-2,\gamma/2}$ be such that $\mathbb{P}\{|\mathcal{T}_{n-2}| > T_{n-2,\gamma/2}\} = \gamma/2$.

- ① A confidence interval of level $1 - \gamma$, for B is

$$\beta \pm T_{n-2,\gamma/2} \times \frac{s_y}{\sqrt{SSX}}.$$

- ② A confidence interval of level $1 - \gamma$, for A is

$$\alpha \pm T_{n-2,\gamma/2} \times s_y \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{SSX}}.$$

- ③ The prediction interval for a new observation x_0 , of level $1 - \gamma$, is

$$\alpha + \beta x_0 \pm T_{n-2,\gamma/2} \times s_y \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SSX}}.$$

Student Distribution

Definition

Let T be a continuous random variable, it follows a t-student distribution with n degree of freedom if its density is given by

$$f(x) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{n\pi}\Gamma\left(\frac{n}{2}\right)} \left(1 + \frac{t^2}{n}\right)^{-\frac{n+1}{2}}.$$

We denote $T \sim \mathcal{T}_n$.

Proposition

Let $Z \sim N(0, 1)$ and $D \sim \chi_n^2$, suppose more that Z and D are independent. Then

$$T = \frac{Z}{\sqrt{D/n}} \sim \mathcal{T}_n.$$

Checking Gaussian hypothesis

To valid our results we need to check our Gaussian assumption about the residual:

- 1 We need to test if the residue comes from a normal distribution, for this we may use the QQ-plot graphical test.
- 2 We need to test if the residue and the fitted come from independent variables, for this we may use the scatterplot of the residue versus the fitted values as a graphical test.

Non-Gaussian case

- All results above are based on the assumption of normality of the residue.
- One may ask whether we still have the same results if the normality assumption is not any more satisfied?
- This leads to the mathematical investigation of looking at the limit behavior of the distribution when the number of observations goes to infinity.
- The results are still approximately correct !
- What we should do when the model is non-gaussian? A suggestion is:
 - 1 add predictor variable(s).
 - 2 and look for non-linear model.
- This will be the topic of the next chapters.