

# Linear regression

It models the relation between "causes" (or independent variable) and "effects" (or dependent variable).

$$(x, y_1, \dots, y_k)' = (\text{effect}, \text{cause 1}, \dots, \text{cause k})$$

for example: ("Sale Price", "Sq Meter Living", "Bathrooms",  
"Bedrooms", "Year Built", "zone")

effect

causes

# Linear regression

The mathematical tool behind this model is the conditional expectation:

$$E(X|Y=y) : \mathbb{R}^k \rightarrow \mathbb{R}, \quad y = (y_1, \dots, y_k)'$$

$$E(X|Y=y) = \int_{\mathbb{R}} x \cdot \frac{f(x, y_1, \dots, y_k)}{f(y_1, \dots, y_k)} dx,$$

where the marginal density of  $Y$  is  $f(y_1, \dots, y_k) = \int_{\mathbb{R}} f(x, y_1, \dots, y_k) dx$ .

The function  $y \rightarrow E(X|Y=y)$  is called the regression of  $X$  in  $Y$ .

# Linear regression

The model of linear regression is the following:

$$E(X|Y=y) = \beta_0 + \sum_{j=1}^k y_j \beta_j,$$

where  $\beta_0, \beta_1, \dots, \beta_k$  are the regression (real) parameters.  
 $\beta_0$  called intercept

The regression hyperplane has the equation:

$$x = \beta_0 + \beta_1 y_1 + \dots + \beta_k y_k$$

Suppose that the statistical data are:

$$(x_i, y_{i1}, \dots, y_{ik}), \quad i = 1, \dots, n \quad n > k$$

# Linear regression

Under the assumption of linearity, the regression parameters can be determined by the Least Squares method:

$$SS(\beta_0, \beta_1, \dots, \beta_k) = \sum_{i=1}^n (x_i - \beta_0 - \beta_1 y_{i1} - \dots - \beta_k y_{ik})^2$$

$$\frac{\partial SS}{\partial \beta_0} = \frac{\partial SS}{\partial \beta_1} = \dots = \frac{\partial SS}{\partial \beta_k} = 0 \quad \Leftrightarrow$$



# Linear regression

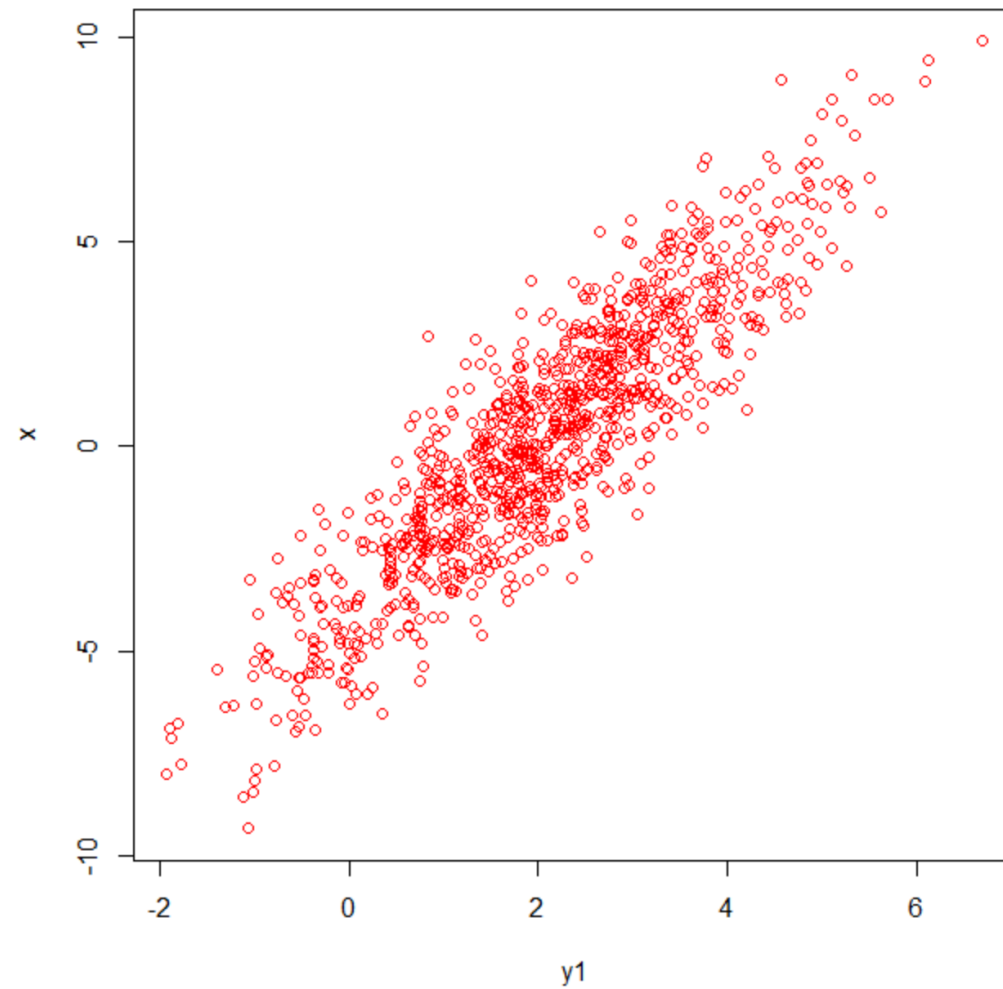
$$\left\{ \begin{array}{l} \sum_{i=1}^n (x_i - \beta_0 - \beta_1 y_{i1} - \dots - \beta_k y_{ik}) = 0 \\ \sum_{i=1}^n y_{ij} (x_i - \beta_0 - \beta_1 y_{i1} - \dots - \beta_k y_{ik}) = 0, \quad j=1, \dots, k \end{array} \right. \Rightarrow$$

$\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$  the LS estimators of the regression parameters  
(the solution is unique if  $\text{rank}(\|y_{ij}\|_{\substack{i=1, \dots, n \\ j=1, \dots, k}}) = k < n$ ).

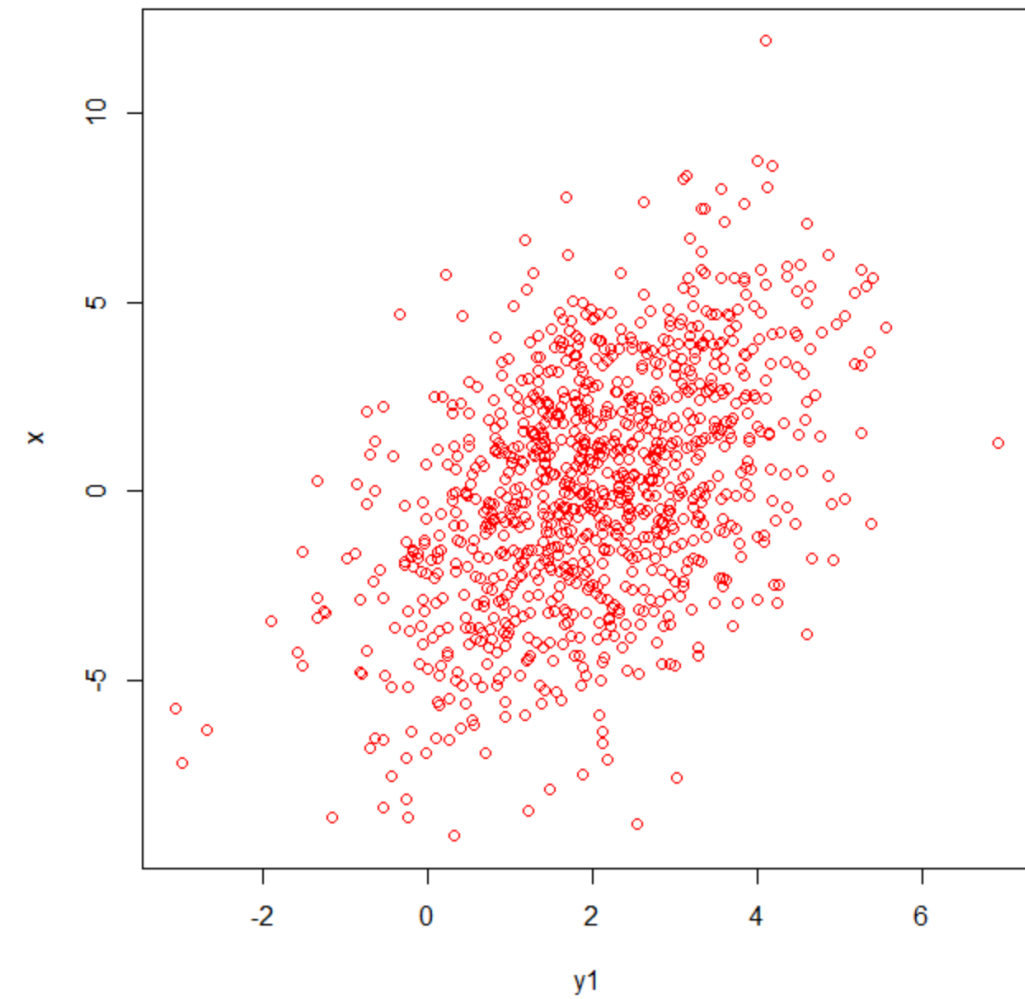
Important obs. This hypothesis of linearity cannot be intuitively identified by visual inspection of the data for  $k \geq 2$ .  
In the case of multivariate normal distribution of the data, the hypothesis of linearity holds.

# Linear regression

```
library(MASS)
mu=c(0,2)
Sigma=matrix(c(10,4,4,2),2,2)
X<-mvrnorm(1000,mu,Sigma)
x<-X[,1]
y1<-X[,2]
plot(y1,x,col="red")
```



```
Sigma=matrix(c(10,2,2,2),2,2)
```



# Linear regression

Proposition Suppose that  $(X, Y_1, \dots, Y_k)'$  is a normally distributed random vector of dimension  $k+1$ :

$$(X, Y_1, \dots, Y_k)' \sim N(\mu, \Sigma), \text{ where}$$

$$\mu = (\mu_x, \mu_y')' \text{ and } \Sigma = \begin{pmatrix} \sigma_x^2 & \Sigma_{x,y} \\ \Sigma_{x,y}' & \Sigma_y \end{pmatrix}, \text{ with } \Sigma_y = \text{Cov}(Y, Y),$$
$$\Sigma_{x,y} = \text{Cov}(X, Y), \quad Y = (Y_1, \dots, Y_k)'$$

Then, the conditional expectation  $E(X|Y)$  is linear:

$$E(X|Y=y) = \Sigma_{x,y} \Sigma_y^{-1} y + (\mu_x - \Sigma_{x,y} \Sigma_y^{-1} \mu_y)$$

The regression hyperplane of  $X$  in  $Y$  can be written as:

$$x - \mu_x = \Sigma_{x,y} \Sigma_y^{-1} (y - \mu_y)$$

# Linear regression

Obs. For  $k=1$ ,  $(X,Y)'$  is bidimensional.

$$(X,Y)' \sim N(\mu, \Sigma), \text{ where } \mu = (\mu_x, \mu_y)', \Sigma = \begin{pmatrix} \sigma_x^2 & \rho \sigma_x \sigma_y \\ \rho \sigma_x \sigma_y & \sigma_y^2 \end{pmatrix}$$

and  $\rho = \frac{\sigma_{x,y}}{\sigma_x \cdot \sigma_y}$  is the correlation coefficient. The conditional

expectation becomes  $E(X|Y=y) = \mu_x + \rho \frac{\sigma_x}{\sigma_y} (y - \mu_y)$  and the

regression line can be written as:

$$x - \mu_x = \rho \cdot \frac{\sigma_x}{\sigma_y} (y - \mu_y).$$



# Linear regression

## Statistical inference for the regression line

For the bidimensional random vector  $(X, Y)'$ , we consider the linear regression model

$$E(X|Y=y) = \beta_0 + \beta_1 y$$

with the regression line

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

Suppose that  $(X_i, Y_i)'_{i=1, \dots, n}$  are observations iid  $\sim (X, Y)'$  and  $(x_i, y_i)'_{i=1, \dots, n}$  are the corresponding statistical data.

# Linear regression

For the statistical data, we consider the linear models

$X_i = \beta_0 + \beta_1 y_i + \xi_i$ ,  $i = 1, \dots, n$ , where  $\xi_1, \dots, \xi_n$  are i.i.d with the conditional mean 0 and the conditional variance  $\sigma^2$ .

We use the Least Square method to estimate the parameters  $\beta_0$  and  $\beta_1$ .

$$SS(\beta_0, \beta_1) = \sum_{i=1}^n (x_i - \beta_0 - \beta_1 y_i)^2$$

$$\frac{\partial SS}{\partial \beta_0} = \frac{\partial SS}{\partial \beta_1} = 0 \Leftrightarrow \begin{cases} n\beta_0 + \beta_1 \sum_{i=1}^n y_i = \sum_{i=1}^n x_i \\ \beta_0 \sum_{i=1}^n y_i + \beta_1 \sum_{i=1}^n y_i^2 = \sum_{i=1}^n x_i y_i \end{cases}$$

The system has a unique solution  $(\hat{\beta}_0, \hat{\beta}_1)$ , assuming that  $\text{not}(y_i = \bar{y}, \forall i)$ .

# Linear regression

Notations :

$$L_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x}) \cdot y_i$$

$$L_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (x_i - \bar{x}) \cdot x_i$$

$$L_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \bar{y}) \cdot y_i$$

The unique solution is:

$$\hat{\beta}_1 = \frac{L_{xy}}{L_{yy}}, \quad \hat{\beta}_0 = \bar{x} - \hat{\beta}_1 \bar{y},$$

and the regression line is:

$$x - \bar{x} = \frac{L_{xy}}{L_{yy}} \cdot (y - \bar{y})$$

# Linear regression

The LS estimators of the parameters of the regression line are:

$$\hat{\beta}_1(x_1, \dots, x_n) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{L_{yy}}$$

$$\hat{\beta}_0(x_1, \dots, x_n) = \bar{x} - \hat{\beta}_1(x_1, \dots, x_n) \bar{y}$$

The estimators have the following properties (in terms of the conditional distribution):

$$E(\hat{\beta}_1) = \beta_1, \quad E(\hat{\beta}_0) = \beta_0,$$

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{L_{yy}}, \quad \text{Var}(\hat{\beta}_0) = \frac{\sigma^2}{n} + \frac{\sigma^2 \bar{y}^2}{L_{yy}},$$

$$\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = -\frac{\sigma^2 \bar{y}}{L_{yy}}.$$



# Linear regression

Now, we perform the analysis of variance for the linear regression ( $k=1$ ).

The fitted (or predicted) values for the regression line are defined as:

$$\hat{x}_i = \hat{\beta}_0 + \hat{\beta}_1 y_i, \quad i = 1, \dots, n$$

We define the following sums of squares:

$$SS_t = \sum_{i=1}^n (x_i - \bar{x})^2$$

$$SS_{\text{regression}} = \sum_{i=1}^n (\hat{x}_i - \bar{x})^2 \quad - \text{how far the slope of the regression line is from 0}$$

$$SS_{\text{residual}} = \sum_{i=1}^n (x_i - \hat{x}_i)^2 \quad - \text{how close the sample points are to the regression line}$$

# Linear regression

For a good-fitting model, we want a high  $SS_{\text{regression}}$  and a low  $SS_{\text{residual}}$ .

$$R^2 \stackrel{\text{def}}{=} \frac{SS_{\text{regression}}}{SS_t} \begin{cases} \text{if the model fits perfect then } R^2=1 \\ \text{if } R^2=0 \text{ it means that } y \text{ gives no} \\ \text{information about } x \end{cases}$$

By direct calculation, we get that:

$$SS_t = L_{xx}, \quad SS_{\text{regression}} = \frac{L_{xy}^2}{L_{yy}}, \quad SS_{\text{residual}} = L_{xx} - \frac{L_{xy}^2}{L_{yy}}$$

(hence  $SS_t = SS_{\text{regression}} + SS_{\text{residual}}$ )

# Linear regression

Notation:  $\overline{SS} = \frac{1}{\text{deg. of freedom}} \cdot \frac{SS}{T^2}$

The hypothesis  $H_0: \{\beta_1 = 0\}$  can be tested by using the F-test.  
Proposition Under the assumption of normality for  $(x, y)'$  and if the hypothesis  $H_0: \{\beta_1 = 0\}$  is true, then

$$\frac{\overline{SS}_{\text{regression}}}{\overline{SS}_{\text{residual}}} = \frac{\frac{1}{1} \cdot \frac{1}{T^2} \cdot SS_{\text{regression}}}{\frac{1}{n-2} \cdot \frac{1}{T^2} \cdot SS_{\text{residual}}} \sim F_{1, n-2}$$

(for details see Dumitrescu & Bătaiorescu, p. 223-224).

# Linear regression

The test statistic for

$$H_0: \{\beta_1 = 0\} \text{ against } H_a: \{\beta_1 \neq 0\}$$

is  $\frac{\overline{SS}_{\text{regression}}}{\overline{SS}_{\text{residual}}}$ , and the acceptance region is

$$W_{n;1-\varepsilon} = \left\{ (x_1, y_1, \dots, x_n, y_n) \mid \frac{\overline{SS}_{\text{regression}}}{\overline{SS}_{\text{residual}}} \leq f_{1,n-2;1-\varepsilon} \right\},$$

where  $\varepsilon$  is the significance level and  $f_{1,n-2;1-\varepsilon}$  is the  $(1-\varepsilon)$  quantile of the  $F_{1,n-2}$  distribution.

The hypothesis  $H_0: \{\beta_1 = 0\}$  against  $H_a: \{\beta_1 \neq 0\}$  can also be tested by a T-test.



# Linear regression

Proposition Under the assumption of normality for  $(x, y)'$  and if the hypothesis  $H_0: \{\beta_1 = 0\}$  is true, then

$$\tilde{t} \stackrel{\text{not}}{=} \frac{\hat{\beta}_1}{\sqrt{\frac{1}{n-2} \cdot SS_{\text{residual}} \cdot \frac{1}{L_{yy}}}} \sim t_{n-2}$$

The acceptance region of the T-test at the significance level  $\varepsilon$  is:

$$W'_{n;1-\varepsilon} = \{(x_1, y_1, \dots, x_n, y_n) \mid -t_{n-2;1-\varepsilon} \leq \tilde{t} \leq t_{n-2;1-\varepsilon}\},$$

where  $t_{n-2;1-\varepsilon}$  is the  $(1-\varepsilon)$  quantile of the  $t_{n-2}$  distribution.

# Linear regression

Interval estimation for  $\beta_0, \beta_1$

Proposition Under the assumption of normality for  $(X, Y)$ , the following random variables are distributed  $t_{n-2}$ :

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{1}{n-2} \cdot SS_{\text{residual}} \cdot \frac{1}{L_{yy}}}} \sim t_{n-2}$$

$$\frac{\hat{\beta}_0 - \beta_0}{\sqrt{\frac{1}{n-2} \cdot SS_{\text{residual}} \cdot \left( \frac{1}{n} + \frac{\bar{y}^2}{L_{yy}} \right)}} \sim t_{n-2}$$

# Linear regression

For the confidence level  $1-\varepsilon$ , the confidence intervals are:

$$C_{n;1-\varepsilon}^{(\beta_1)} = \left\{ -t_{n-2;1-\varepsilon} \leq \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{1}{n-2} SS_{\text{residual}} \cdot \frac{1}{L_{yy}}}} \leq t_{n-2;1-\varepsilon} \right\}$$

$$= \left\{ \hat{\beta}_1 - t_{n-2;1-\varepsilon} \cdot \sqrt{\frac{1}{n-2} \cdot SS_{\text{residual}} \cdot \frac{1}{L_{yy}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{n-2;1-\varepsilon} \cdot \sqrt{\frac{1}{n-2} \cdot SS_{\text{residual}} \cdot \frac{1}{L_{yy}}} \right\}$$

$$C_{n;1-\varepsilon}^{(\beta_0)} = \left\{ -t_{n-2;1-\varepsilon} \leq \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\frac{1}{n-2} \cdot SS_{\text{residual}} \cdot \left( \frac{1}{n} + \frac{\bar{y}^2}{L_{yy}} \right)}} \leq t_{n-2;1-\varepsilon} \right\}$$

$$= \left\{ \dots \leq \beta_0 \leq \dots \right\}$$

# Linear regression

The sample correlation coefficient  $r = \frac{L_{xy}}{\sqrt{L_{xx}L_{yy}}}$ .

The relationship between  $r$  and  $\hat{\beta}_1$  is:

$$\hat{\beta}_1 = r \cdot \frac{\sqrt{L_{xx}}}{\sqrt{L_{yy}}}$$

We have that  $-1 \leq r \leq 1$ .

$\left\{ \begin{array}{l} r > 0 \text{ (} r < 0 \text{)} \text{ then } X \text{ tends to increase (resp. decrease) as } Y \text{ increases} \\ r = 0 \text{ then } X \text{ is unrelated to } Y. \end{array} \right.$

if we are interested in whether there is an association between two variables, then we analyze  $r$ .

if we are interested in prediction, then we take a look at  $\hat{\beta}_1$ .



# Linear regression

## Summary (for $k=1$ )

- we have a set of data  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ ,  $x_i, y_i \in \mathbb{R}$
- we fit a linear regression model to our data - we estimate the parameters  $\beta_0$  and  $\beta_1$
- if we consider that  $SS_{\text{residual}}$  is "small enough", then we make the assumption that the data are normally distributed.

# Linear regression

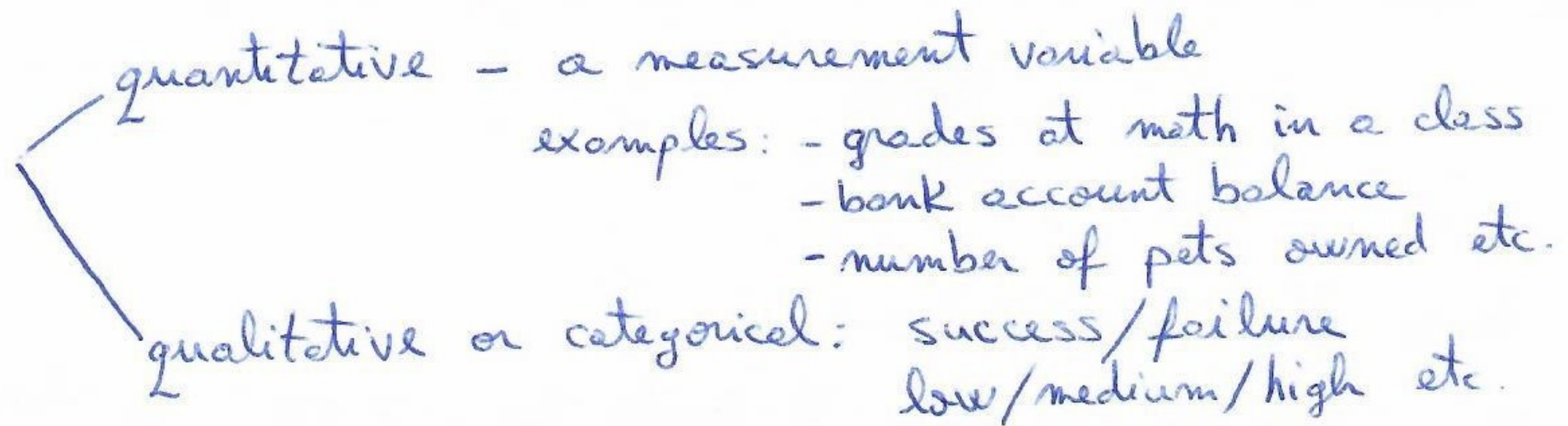
- under the assumption of normality, we perform ANOVA.

This analysis is important because if  $\beta_1 = 0$  then our data are uncorrelated — either the linearity hypothesis does not hold or the data are independent — in both cases, we will not apply a linear regression model for our data.

# Linear regression

## Comment

There are two types of variables used in statistics:



Regression = model for the relationship "causes" → "effects"



# Linear regression

Depending on the type of the causes and effects, we have different situations that are differently treated:

- 1) The effect is qualitative; the causes are quantitative  
- the typical model used is the Generalized Linear Regression (e.g. Logistic Regression, Log-Linear Regression)
- 2) The effect is qualitative/quantitative; the causes are qualitative - the typical model used is Factor Analysis.
- 3) The effect is quantitative; the causes are quantitative - the typical model used is the Linear Regression.