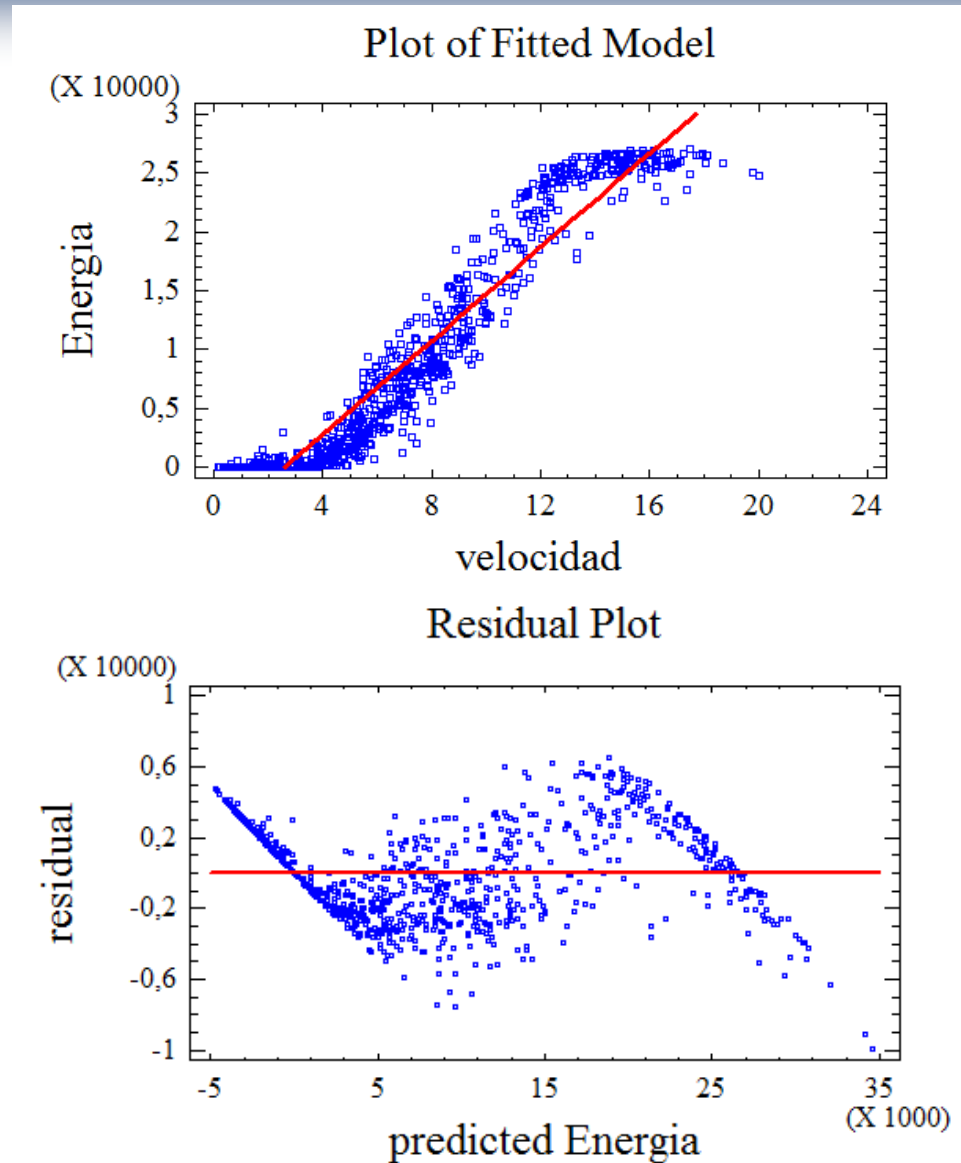


# II. BIVARIATE DESCRIPTIVE STATISTICS



# Chapter II: Bivariate Descriptive Statistics

- 1. Introduction.**
- 2. Bivariate Frequency Tables**
- 3. Scatterplots**
- 4. Measures of linear dependence**
- 5. The regression line**

# Chapter II: Bivariate Descriptive Statistics

1. Introduction.
2. **Bivariate Frequency Tables**
3. Scatterplots
4. Measures of linear dependence
5. The regression line

## 2. Bivariate Frequency Tables

### Bivariate tables

We have, for each individual, two variables and to describe them we use a table with a double entrance

Example: for each car we have the number of cylinders and its manufacturing year (file cardata.sf)

	78	79	80	81	82	Row Total
3	0	0	1	0	0	1
4	17	12	25	22	28	104
5	1	1	1	0	0	3
6	12	6	2	7	3	30
8	6	10	0	1	0	17
Column Total	36	29	29	30	31	155

## 2. Bivariate Frequency Tables

### Bivariate tables

We have, for each individual, two variables and to describe them we use a table with a double entrance

Example: for each car we have the number of cylinders and its manufacturing year (file cardata.sf)

	78	79	80	81	82	Row Total
3	0	0	1	0	0	1
4	17	12	25	22	28	104
5	1	1	1	0	0	3
6	12	6	2	7	3	30
8	6	10	0	1	0	17
Column Total	36	29	29	30	31	155

Each cell contains the **absolute joint frequency**

## 2. Bivariate Frequency Tables

### Bivariate tables

We have, for each individual, two variables and to describe them we use a table with a double entrance

Example: for each car we have the number of cylinders and its manufacturing year (file cardata.sf)

	78	79	80	81	82	Row Total
3	0	0	1	0	0	1
4	17	12	25	22	28	104
5	1	1	1	0	0	3
6	12	6	2	7	3	30
8	6	10	0	1	0	17
Column Total	36	29	29	30	31	155

Univariant: **absolute marginal frequencies**

## 2. Bivariate Frequency Tables

### Bivariate tables

We have, for each individual, two variables and to describe them we use a table with a double entrance

Example: for each car we have the number of cylinders and its manufacturing year (file cardata.sf)

	78	79	80	81	82	Row Total
3	0	0	1	0	0	1
4	17	12	25	22	28	104
5	1	1	1	0	0	3
6	12	6	2	7	3	30
8	6	10	0	1	0	17
Column Total	36	29	29	30	31	155

Each row or column contains the **absolute conditional frequency** (with respect to the value of the row or column)

# Chapter II: Bivariate Descriptive Statistics

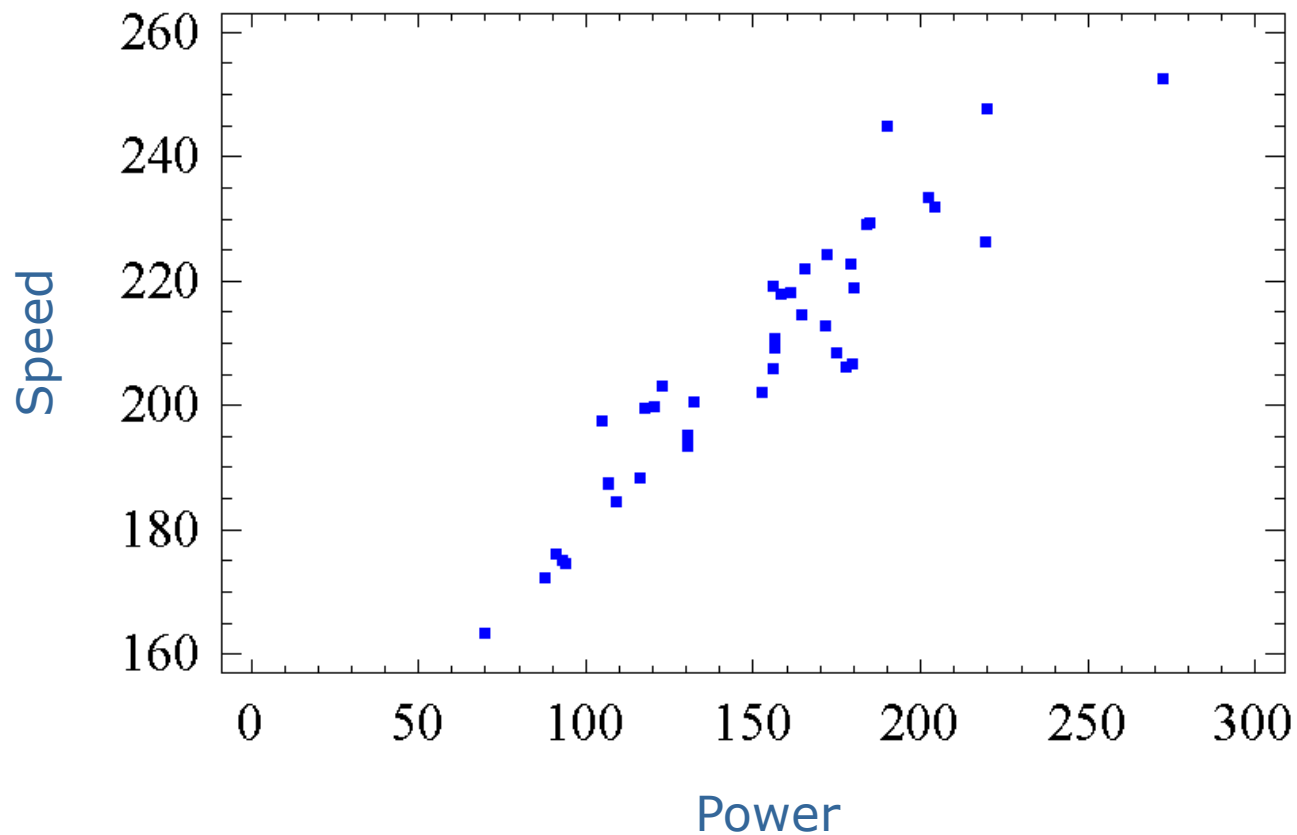
1. Introduction.
2. Bivariate Frequency Tables
3. Scatterplots
4. Measures of linear dependence
5. The regression line



### 3. Scatterplot

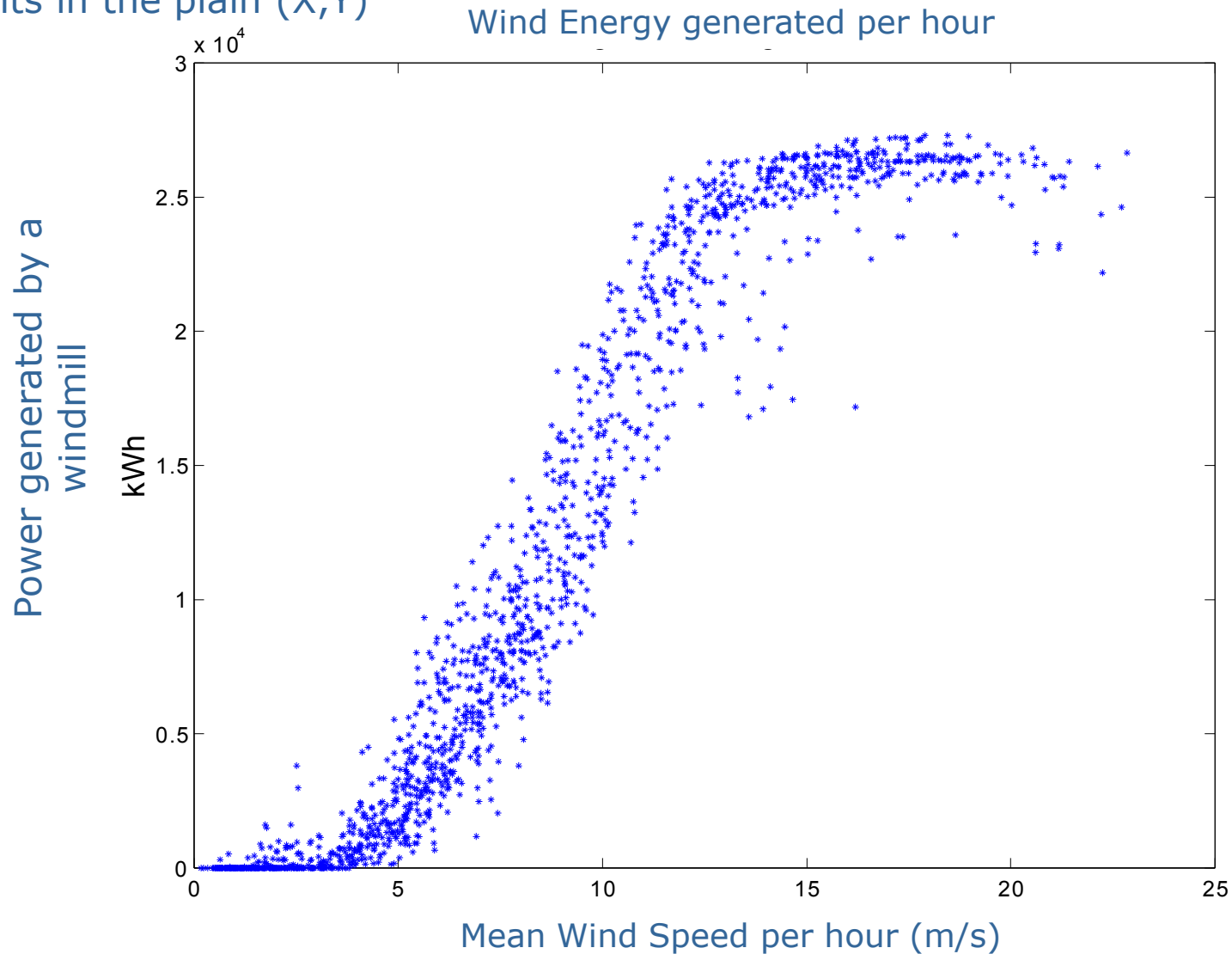
For each individual we have two data, X and Y, and we plot all data as points in the plain (X,Y)

Scatter-plot – Speed vs. Power



### 3. Scatterplot

For each individual we have two data, X and Y, and we plot all data as points in the plane (X,Y)



# Chapter II: Bivariate Descriptive Statistics

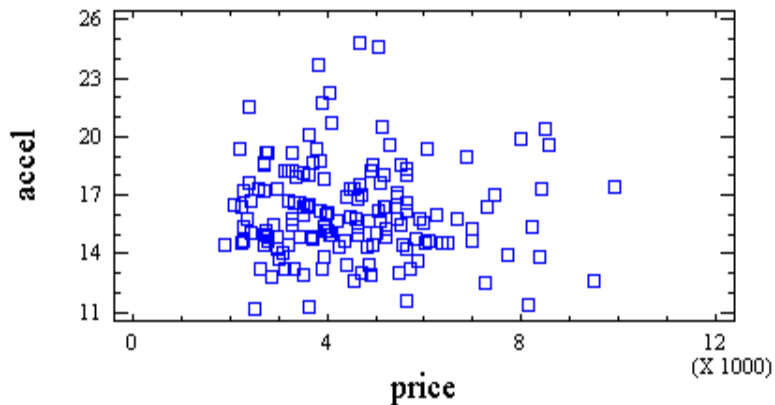
1. Introduction.
2. Bivariate Frequency Tables
3. Scatterplots
4. Measures of linear dependence
5. The regression line

### 3. Measures of linear dependence

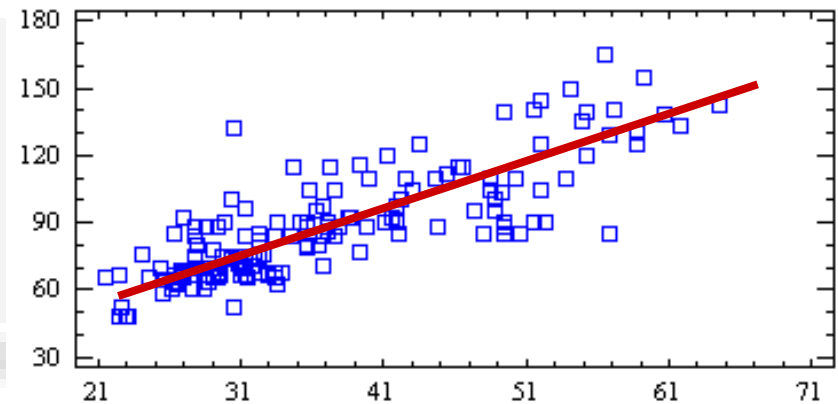
#### Measures of linear dependence

- Covariance coefficient
- Correlation coefficient

Plot of accel vs price



Between these variables  
does not exist a linear  
relation



Between these variables  
exist a linear relation

The red line could be a good  
summary of that relation

For n individuals, we have data of 2 variables

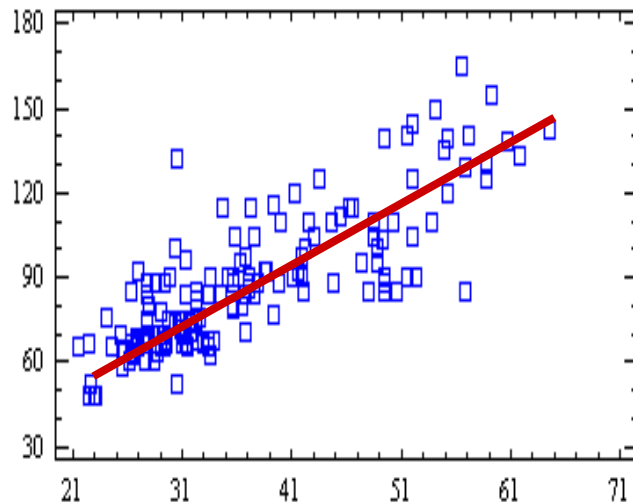
Individuals	<b>x</b>	<b>y</b>
1	x1	y1
2	x2	y2
:	:	:
n	xn	yn

### Covariance

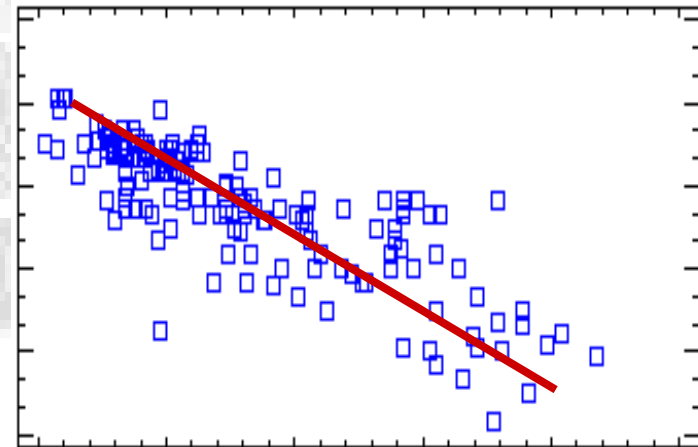
$$\text{cov}(x,y) = s_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n}$$

### Correlation

$$r = r_{xy} = r(x,y) = \frac{\text{cov}(x,y)}{s_x s_y}$$



Covariance and  
positive correlation



Covariance and  
negative correlation

### 3. Measures of linear dependence

A usual way to describe this information in a (symmetric) matricial form is by using the

Covariance Matrix

$$M = \begin{bmatrix} s_x^2 & \text{cov}(x, y) \\ \text{cov}(y, x) & s_y^2 \end{bmatrix}$$

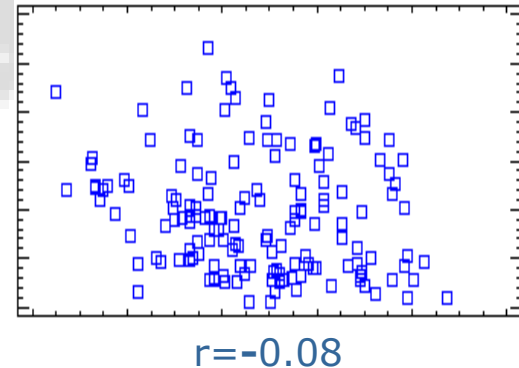
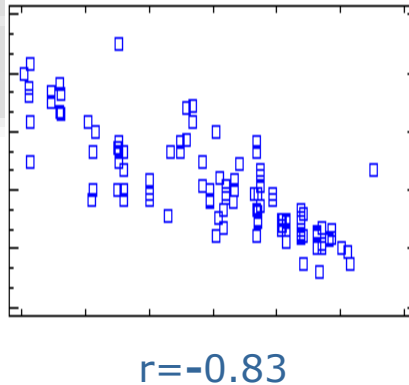
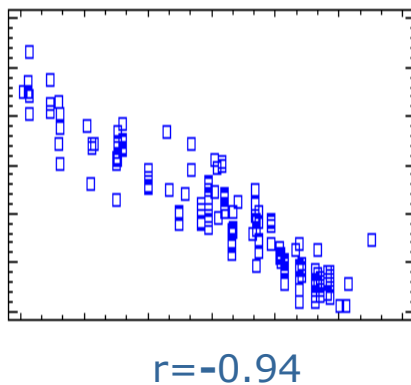
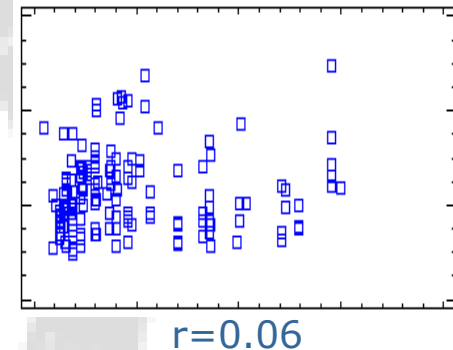
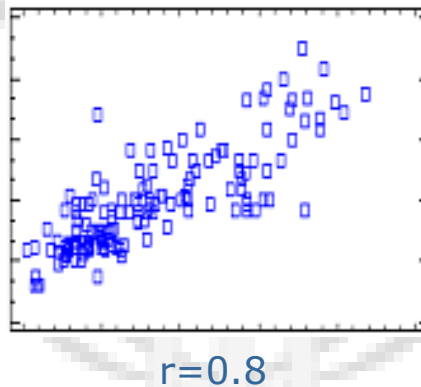
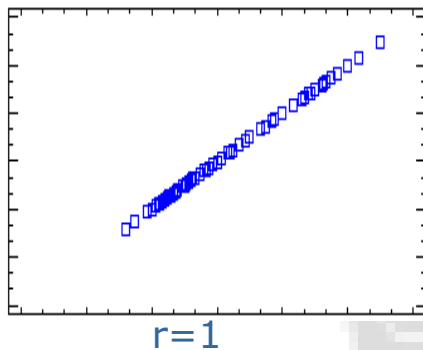
Correlation Matrix

$$R = \begin{bmatrix} 1 & \text{corr}(x, y) \\ \text{corr}(y, x) & 1 \end{bmatrix}$$

$$\text{cov}(x,y) = s_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n}$$

$$r = r_{xy} = r(x,y) = \frac{\text{cov}(x,y)}{s_x s_y}$$

- The covariance depends upon the units in which x and y are measured
- The correlation is dimensionless. IT IS EASIER TO INTERPRET
- The correlation coefficient always is bounded  $-1 \leq r \leq 1$

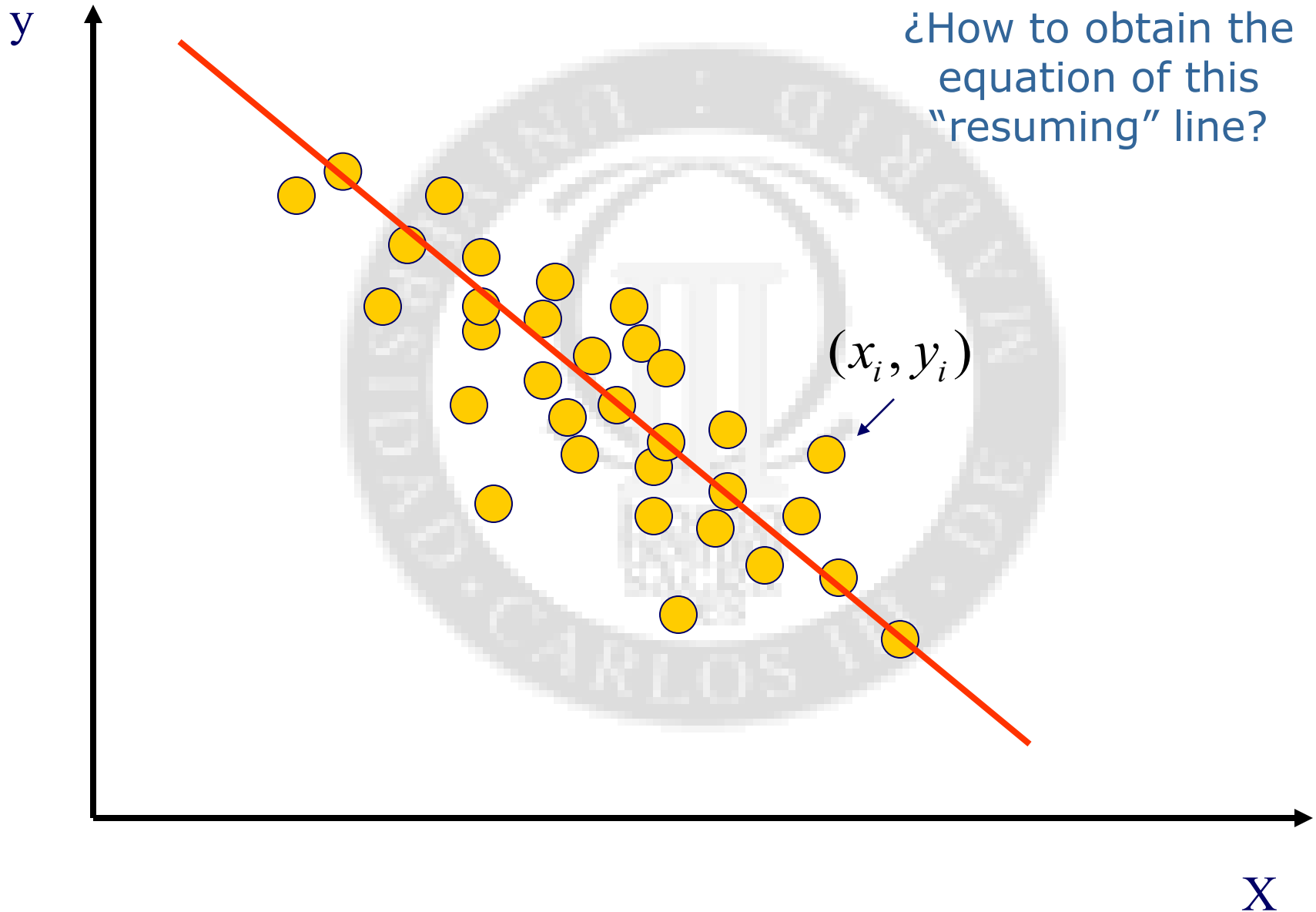


# Chapter II: Bivariate Descriptive Statistics

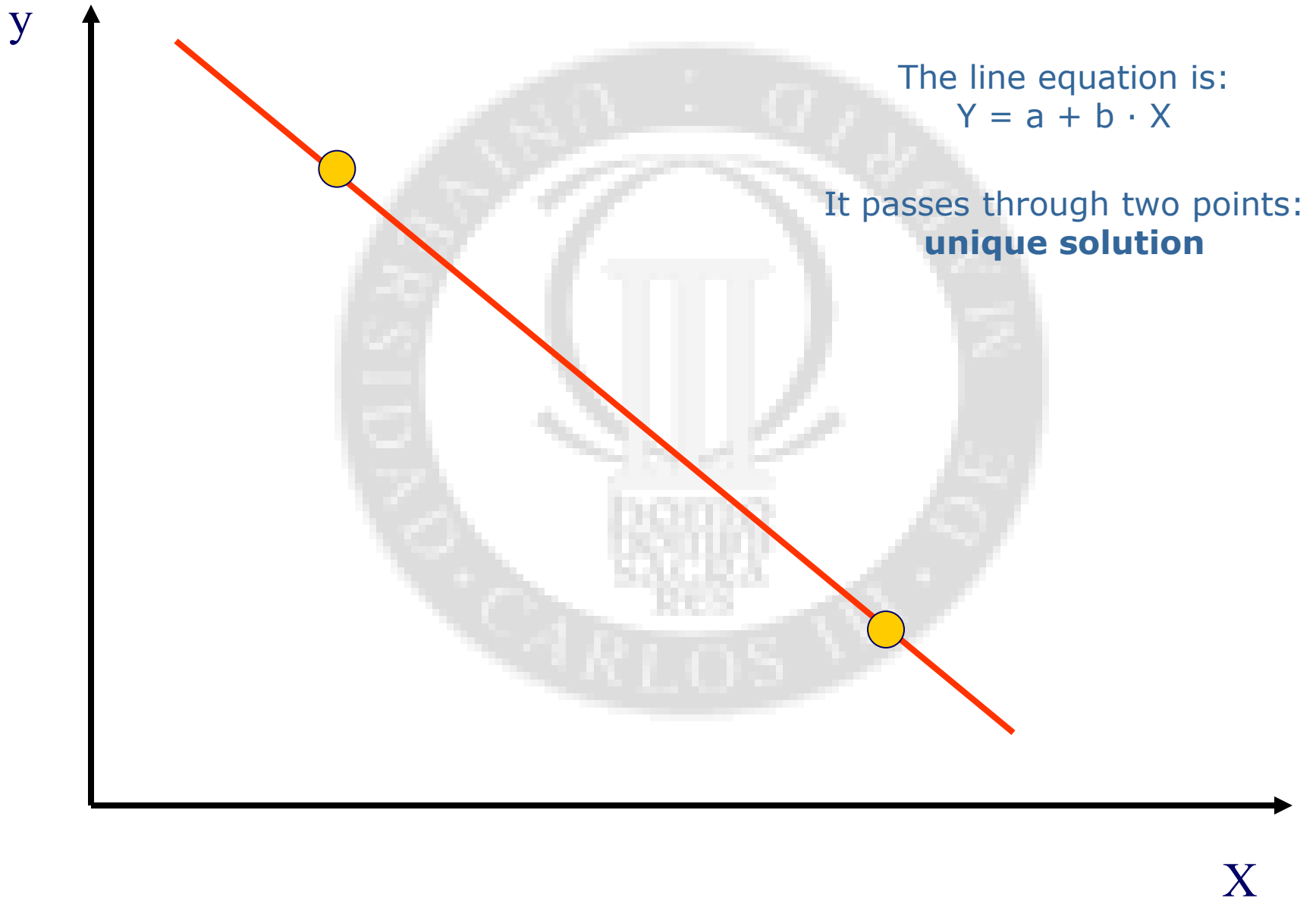
1. Introduction.
2. Bivariate Frequency Tables
3. Scatterplots
4. Measures of linear dependence
5. The regression line



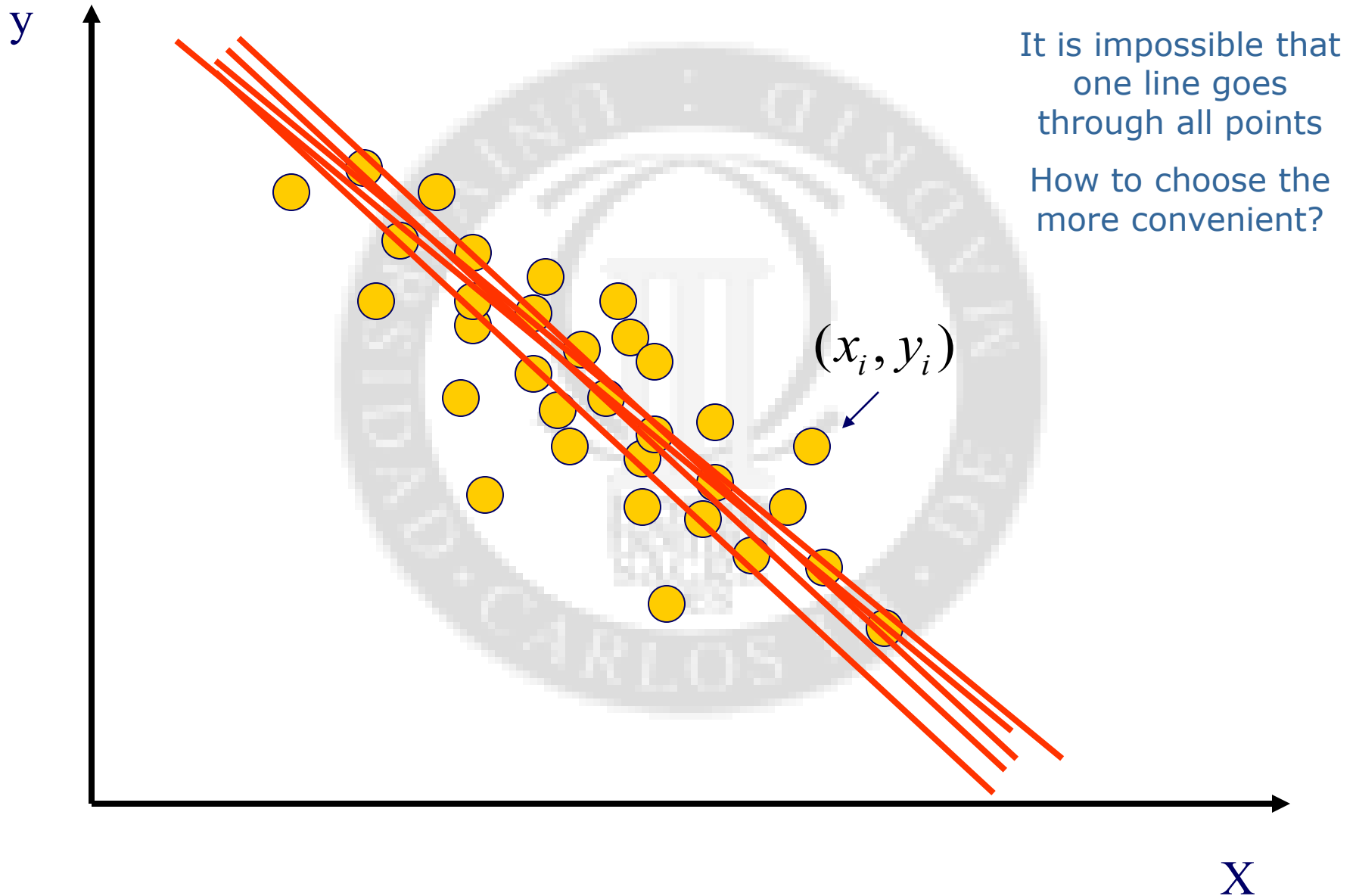
## 4. The regression line



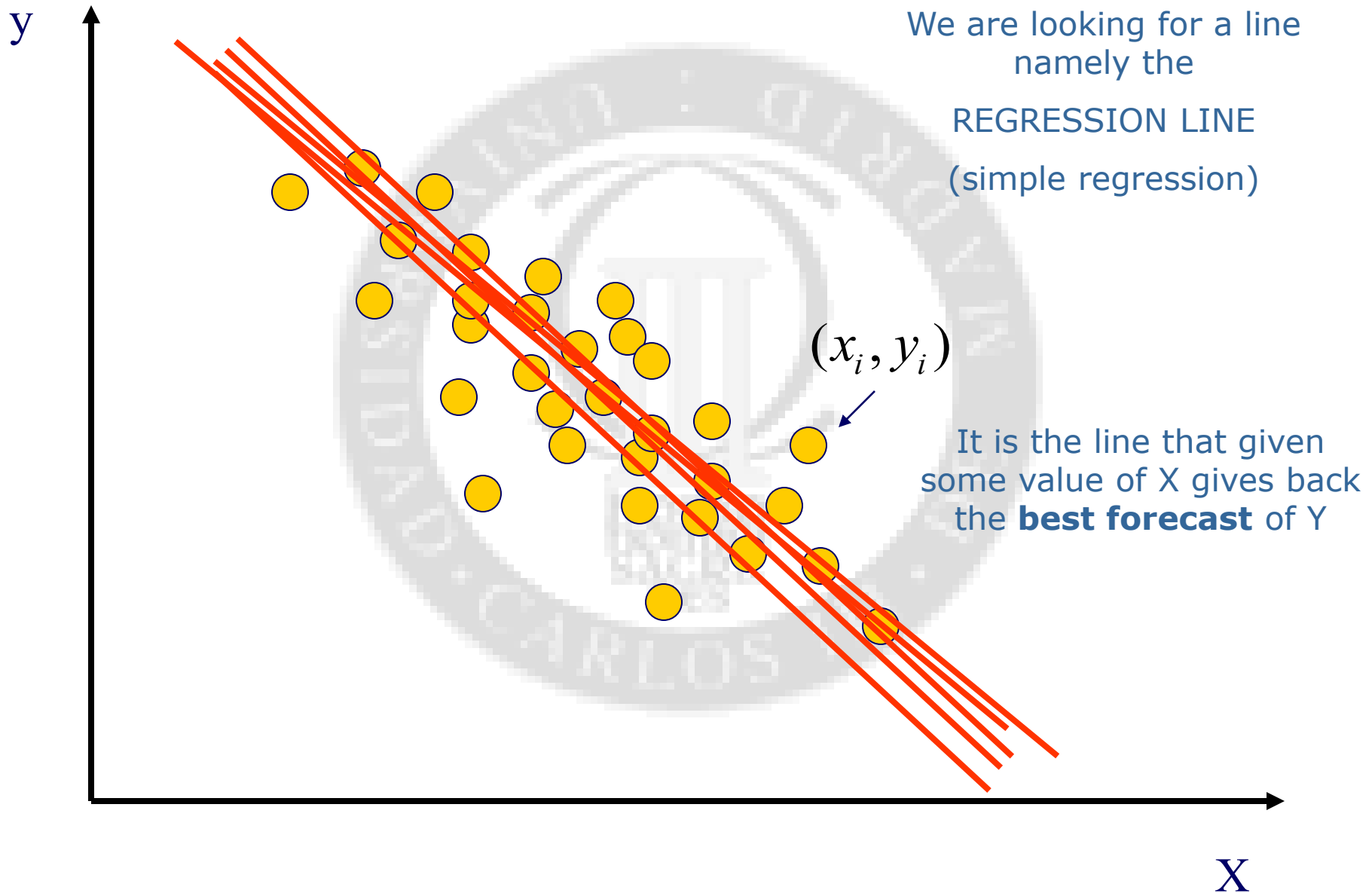
## 4. The regression line



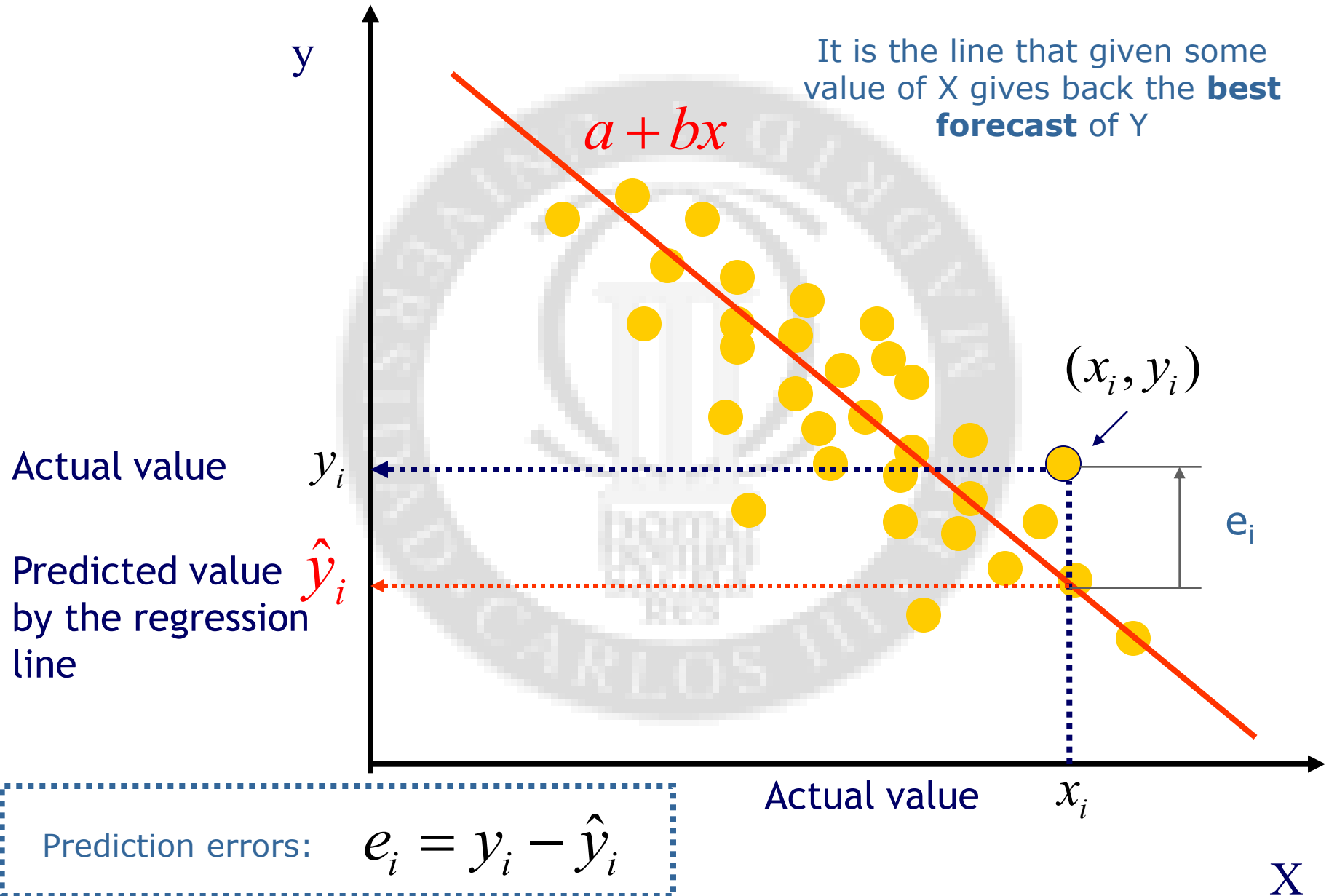
## 4. The regression line



## 4. The regression line



## 4. The regression line

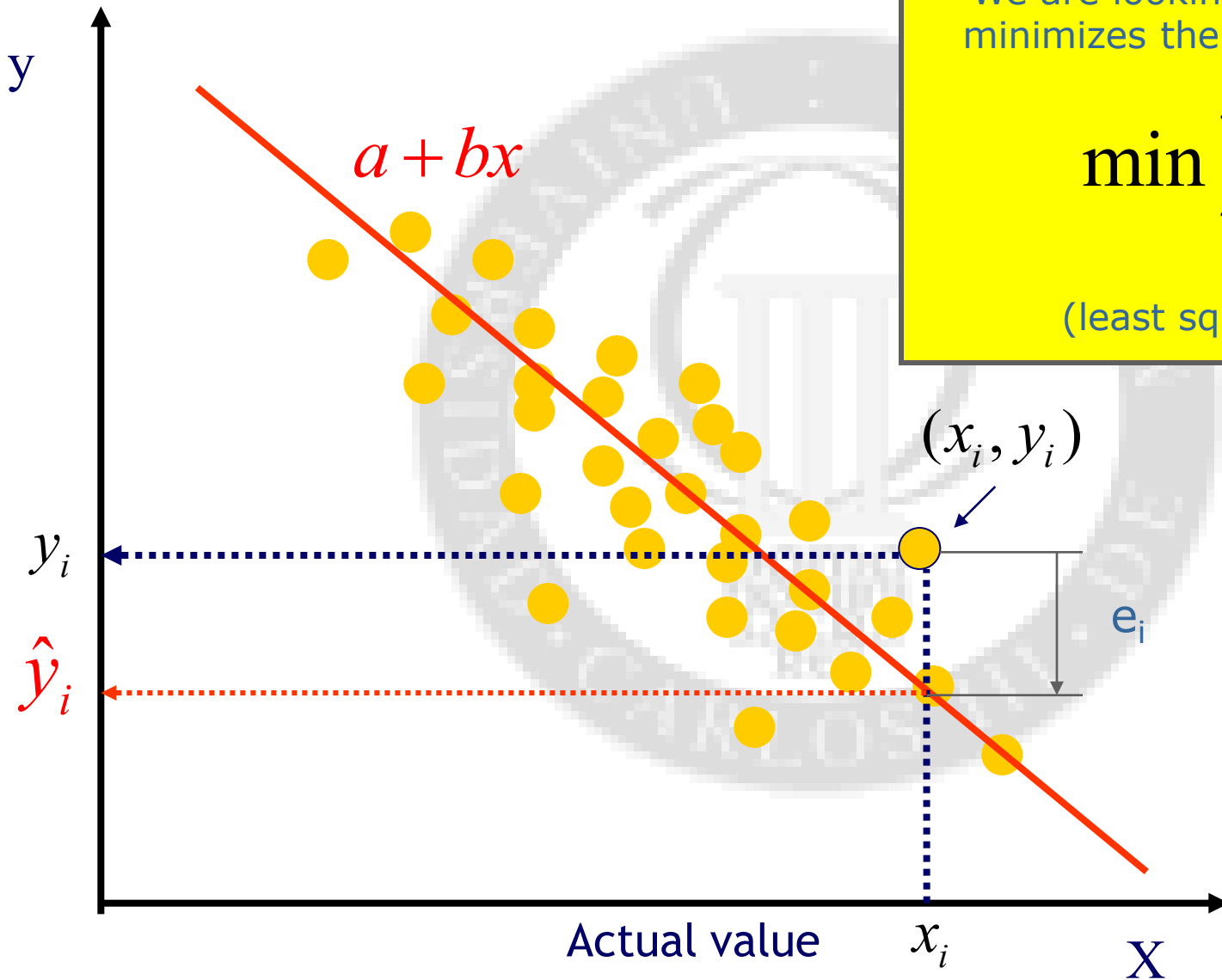


## 4. The regression line

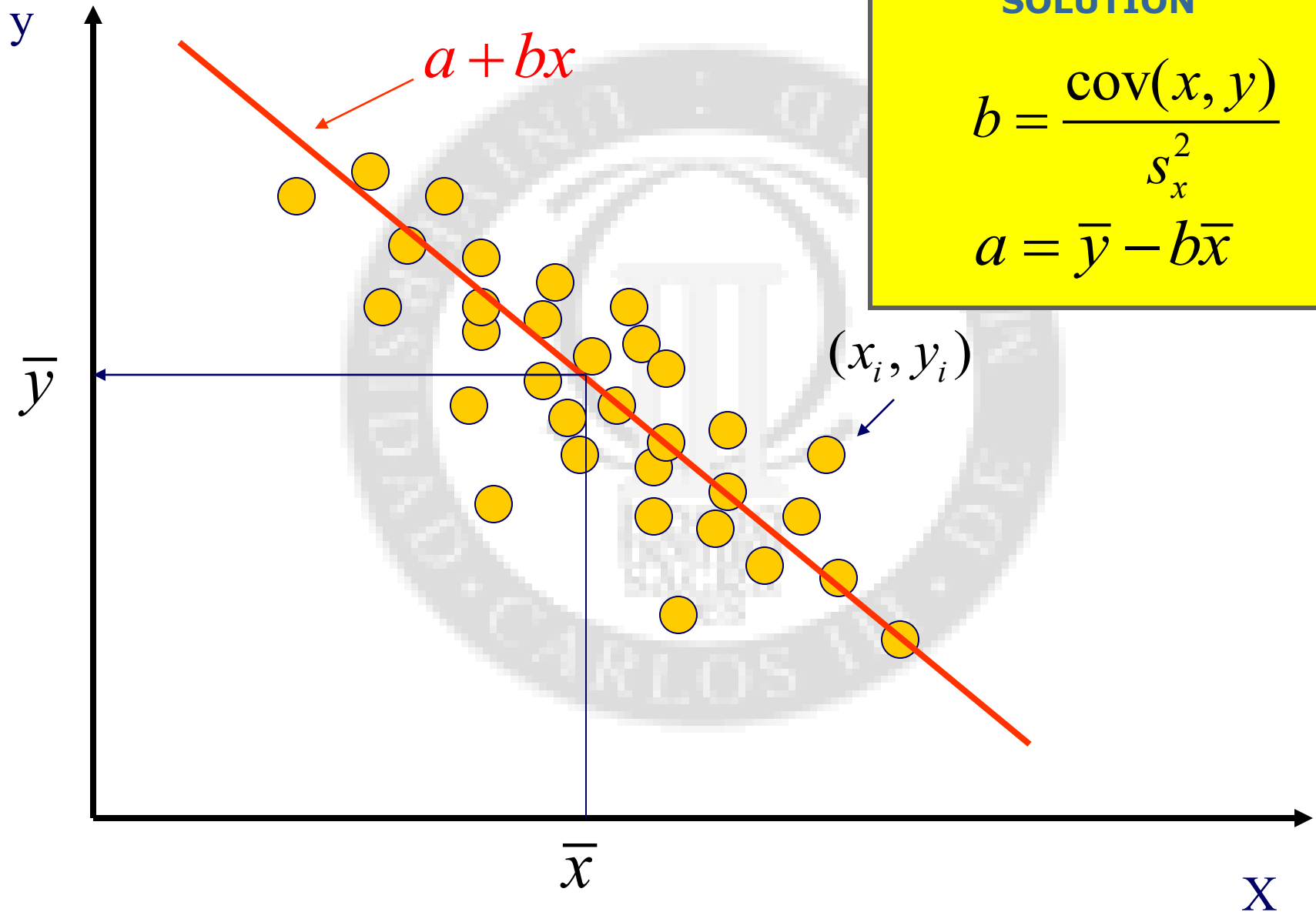
We are looking for the line that minimizes the prediction errors:

$$\min \sum_{i=1}^N e_i^2$$

(least squares line)



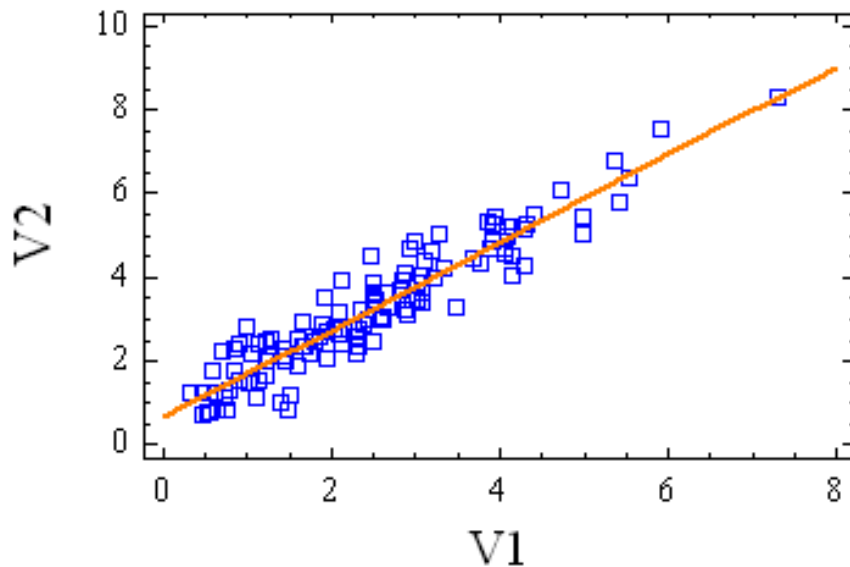
## 4. The regression line



## Example

The variable V1 is the wind speed registered in location 1, while the variable V2 is the speed registered at the same time in location 2. There are a total of 115 pairs of measures.

Plot of Fitted Model



Loc.1:  
mean: 2.51  
variance: 1.91

Loc.2:  
mean: 3.28  
variance: 2.36

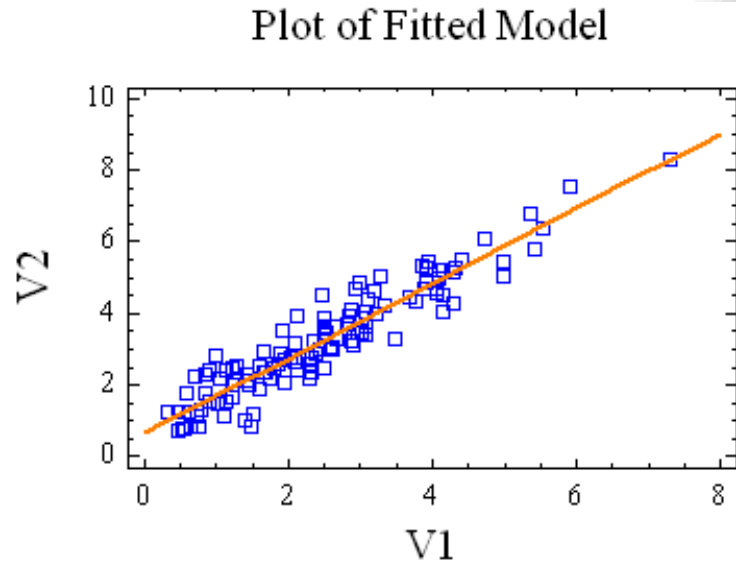
$\text{cov}(V1, V2) = 1.995$

In location 1 it is going to install a computer system to telemeasure the wind speed, but not for location 2. We want to calculate the regression line which allows to predict the speed in location 2 knowing the speed in location 1.



## Example

The variable V1 is the wind speed registered in location 1, while the variable V2 is the speed registered at the same time in location 2. There are a total of 115 pairs of measures.



Loc.1:

mean: 2.51

variance: 1.91

Loc.2:

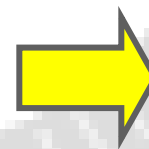
mean: 3.28

variance: 2.36

cov (V1,V2)=1.995

$$b = \text{cov}(x,y) / \text{var}(x) = 1.995 / 1.91 = 1.045$$

$$a = \bar{y} - b\bar{x} = 3.28 - 1.045 \times 2.51 = 0.657$$



$$\hat{V}_2 = 0.657 + 1.045 \times V_1$$

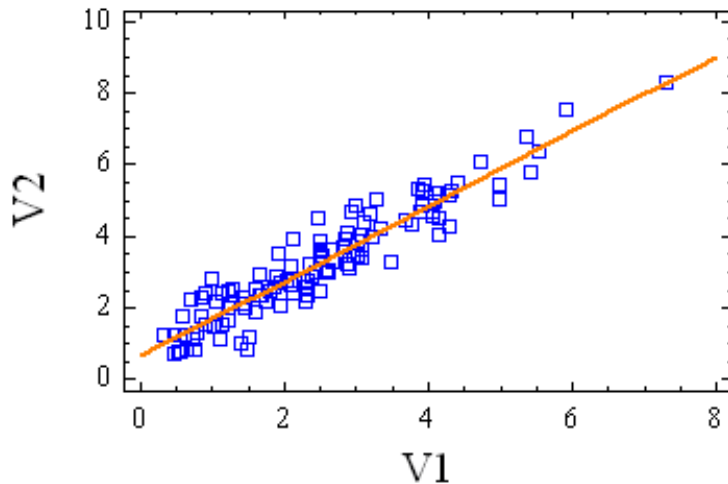
If, for example, in Location 1 we measure a speed wind of 5 m/s, the wind speed prediction for Location 2 is

$$\mathbf{0.657 + 1.045 \times 5 = 5.88 \text{ m/s}}$$

## Example

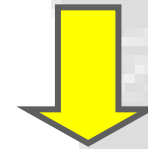
The variable V1 is the wind speed registered in location 1, while the variable V2 is the speed registered at the same time in location 2. There are a total of 115 pairs of measures.

Plot of Fitted Model



$$b = \text{cov}(x, y) / \text{var}(x) = 1.995 / 1.91 = 1.045$$

$$a = \bar{y} - b\bar{x} = 3.28 - 1.045 \times 2.51 = 0.657$$



$$\hat{V}_2 = 0.657 + 1.045 \times V_1$$

$$y = a + bx$$

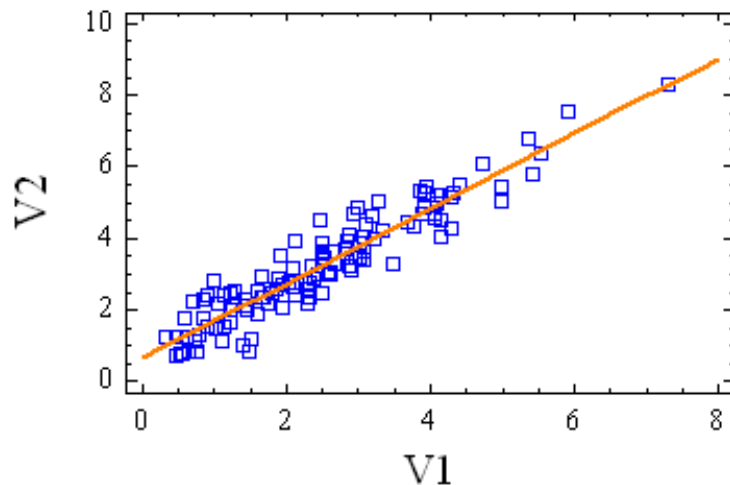
**Interpretation of b:** if x increases of one unit y increases of b units

If in location 1 the wind speed increases of 1 m/s, in location 2 it is predicted to increase of 1.045 m/s

## Example

The variable V1 is the wind speed registered in location 1, while the variable V2 is the speed registered at the same time in location 2. There are a total of 115 pairs of measures.

Plot of Fitted Model



$$b = \text{cov}(x, y) / \text{var}(x) = 1.995 / 1.91 = 1.045$$

$$a = \bar{y} - b\bar{x} = 3.28 - 1.045 \times 2.51 = 0.657$$



$$\hat{V}_2 = 0.657 + 1.045 \times V_1$$

$$y = a + bx$$

**Interpretation of a:** if the value of x is 0, y has value a

If in location 1 there is no wind, in location 2 there would be a wind speed of 0.657 m/s, that is a small value

## 4. The regression line

### Evaluation of the regression

The regression to predict  $y$  starting from the value  $x$  will be good if:

1. The relation between  $x$  and  $y$  is linear
2. The linear relation is strong (big correlation)

Correlation and  
Square of correlation coefficient  $R^2$

Graphs

## Evaluation of the regression

The regression to predict  $y$  starting from the value  $x$  will be good if:

### 1. The relation between $x$ and $y$ is linear

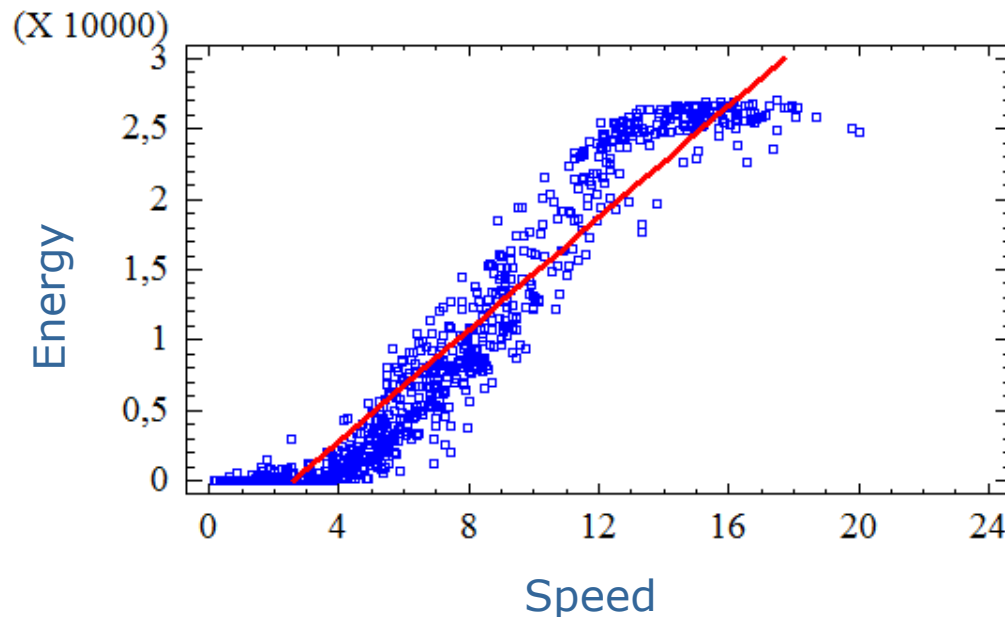
1.1 Scatterplot

1.2 Graph of predictions vs. observation

1.3 Graph of residuals vs. predicted values

#### 1.1 Scatterplot

Plot of Fitted Model



Data: `parqueeolico.sf3`

This simple scatterplot shows that there is no linear relation. The regression line gives bad predictions

# Evaluation of the regression

The regression to predict  $y$  starting from the value  $x$  will be good if:

## 1. The relation between $x$ and $y$ is linear

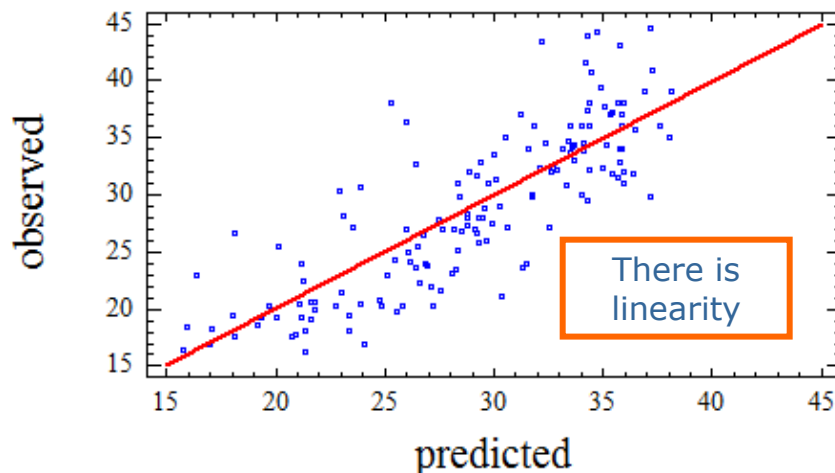
1.1 Scatterplot

1.2 Graph of predictions vs. observation

1.3 Graph of residuals vs. predicted values

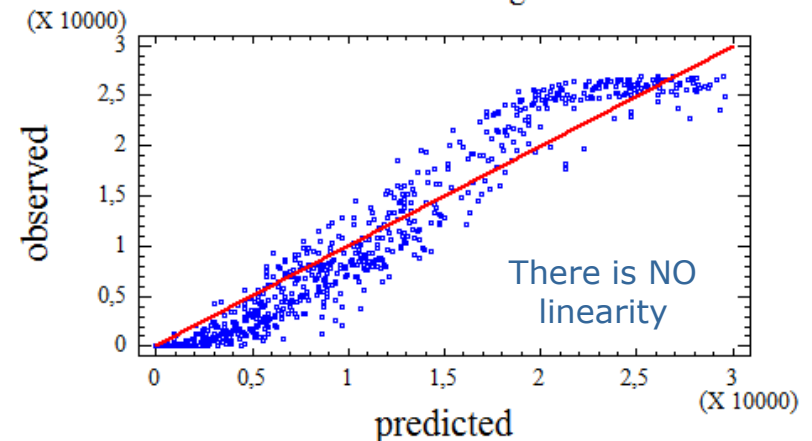
### 1.2 Graph of predictions vs observation

Plot of mpg



Cardata.sf: we want to express mpg (miles per gallon) as function of the weight (weight)

Plot of Energia

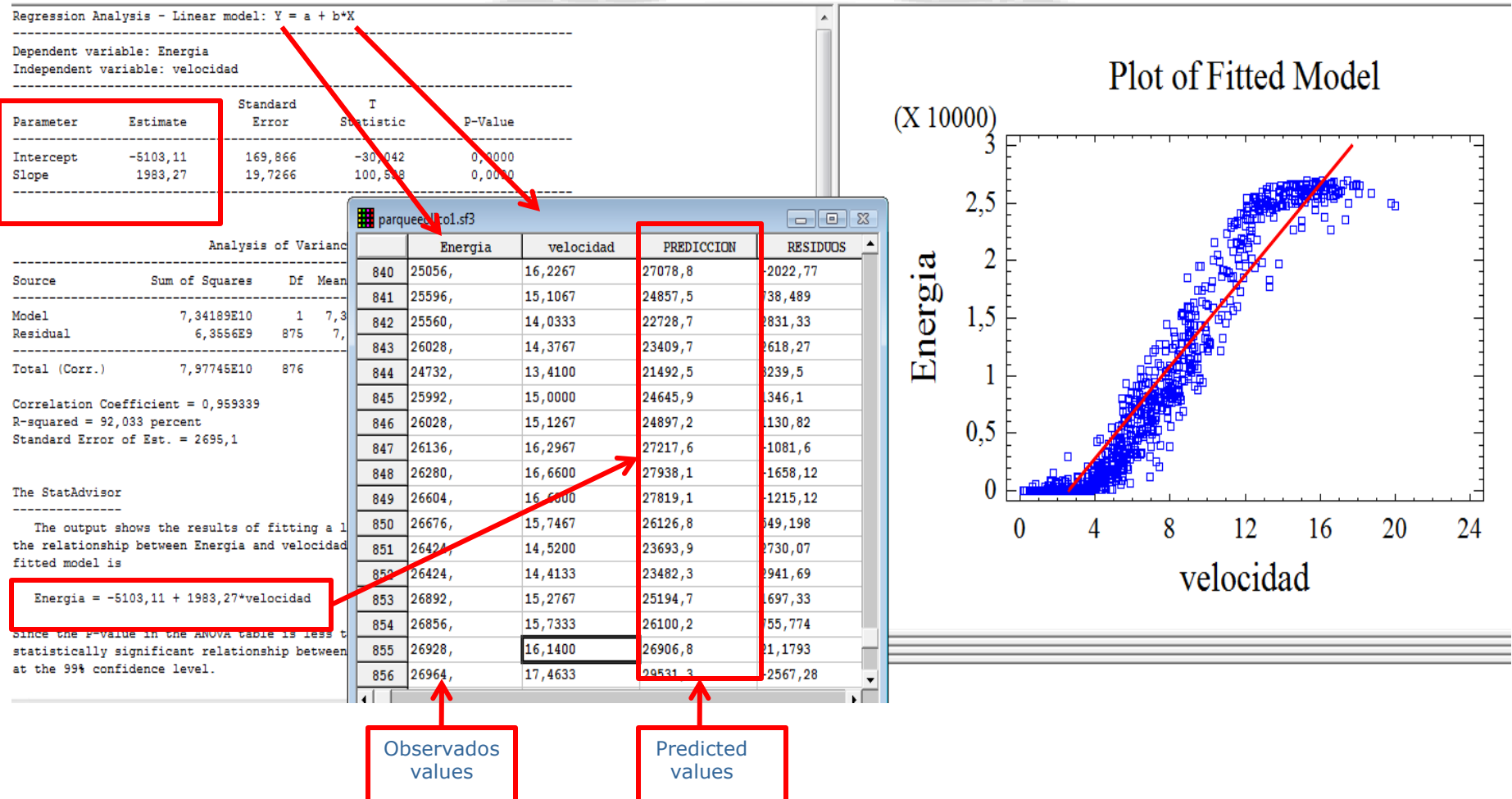


Parqueolico.sf3: we want to express the generated energy as function of the wind speed

# Evaluation of the regression

## 1.3 Graph of residuals vs. predicted values

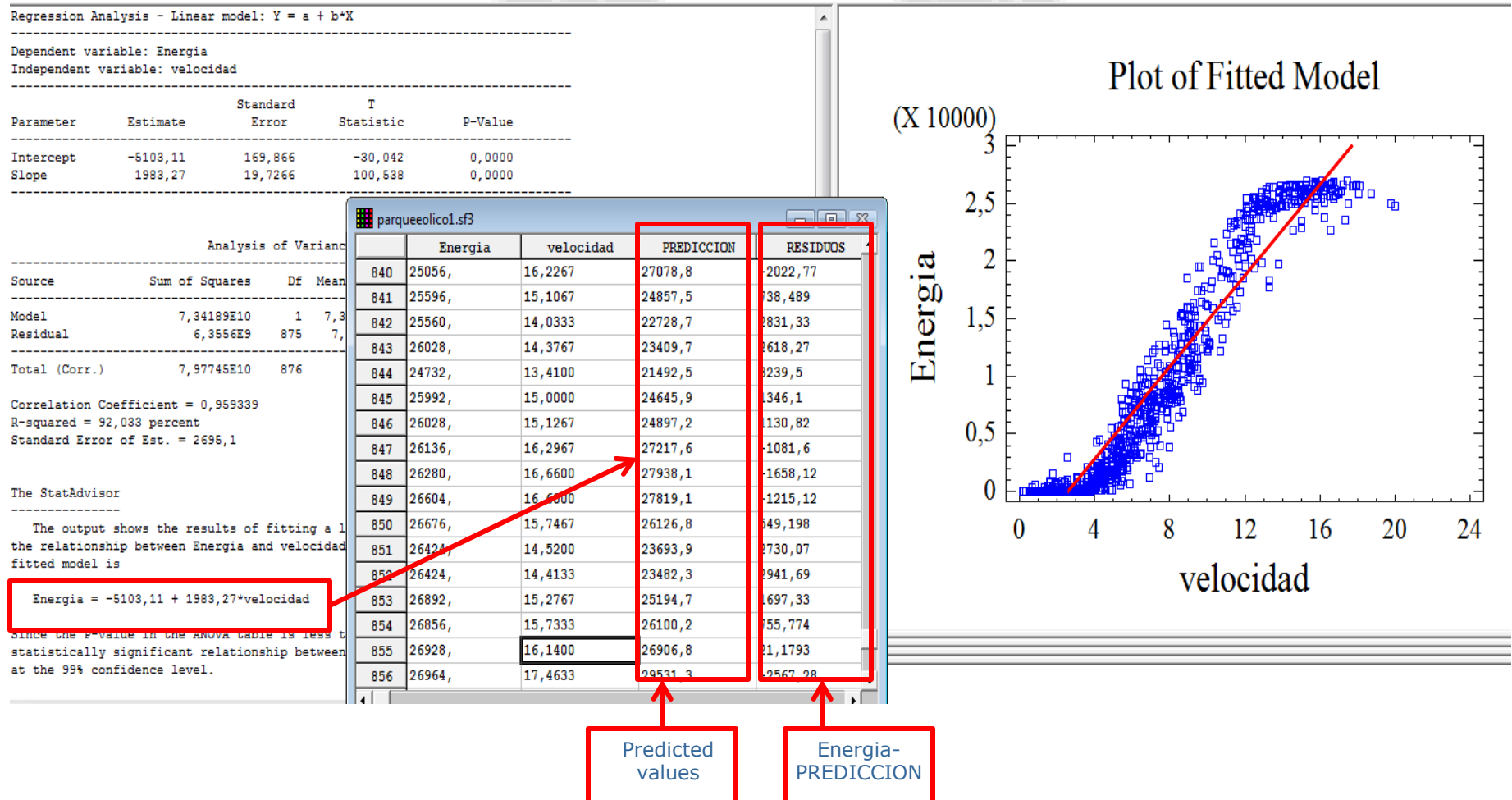
It is the most important graph representation to evaluate a regression



# Evaluation of the regression

## 1.3 Graph of residuals vs. predicted values

It is the most important graph representation to evaluate a regression





b\*X

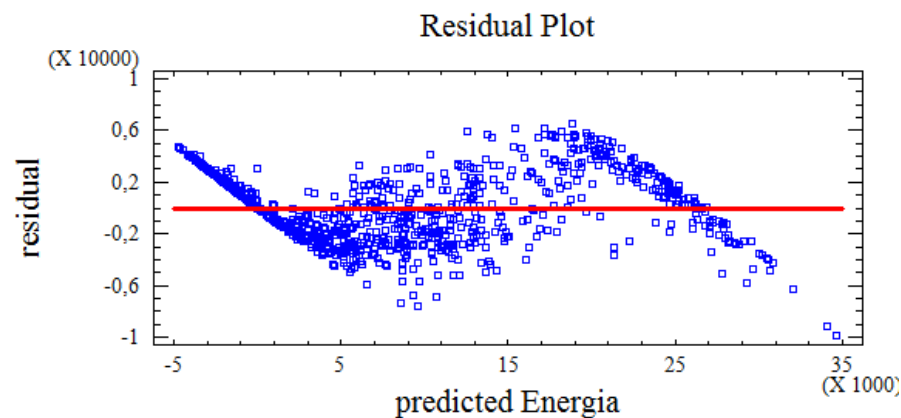
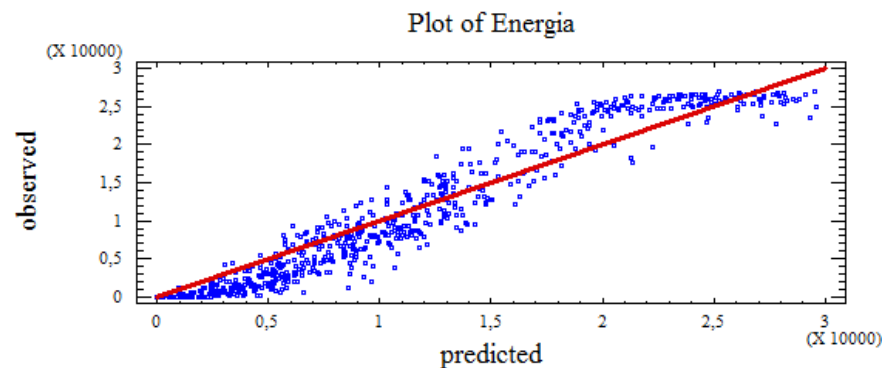
parqueolico1.sf3

	Energia	velocidad	PREDICCION	RESIDUOS
840	25056,	16,2267	27078,8	-2022,77
841	25596,	15,1067	24857,5	738,489
842	25560,	14,0333	22728,7	2831,33
843	26028,	14,3767	23409,7	2618,27
844	24732,	13,4100	21492,5	3239,5
845	25992,	15,0000	24645,9	1346,1
846	26028,	15,1267	24897,2	1130,82
847	26136,	16,2967	27217,6	-1081,6
848	26280,	16,6600	27938,1	-1658,12
849	26604,	16,6000	27819,1	-1215,12
850	26676,	15,7467	26126,8	549,198
851	26424,	14,5200	23693,9	2730,07
852	26424,	14,4133	23482,3	2941,69
853	26892,	15,2767	25194,7	1697,33
854	26856,	15,7333	26100,2	755,774
855	26928,	16,1400	26906,8	21,1793
856	26964,	17,4633	29531,3	-2567,28

s then 0.01, there is a  
een Energia and velocidad

Predicted  
values

Energia-  
PREDICCION

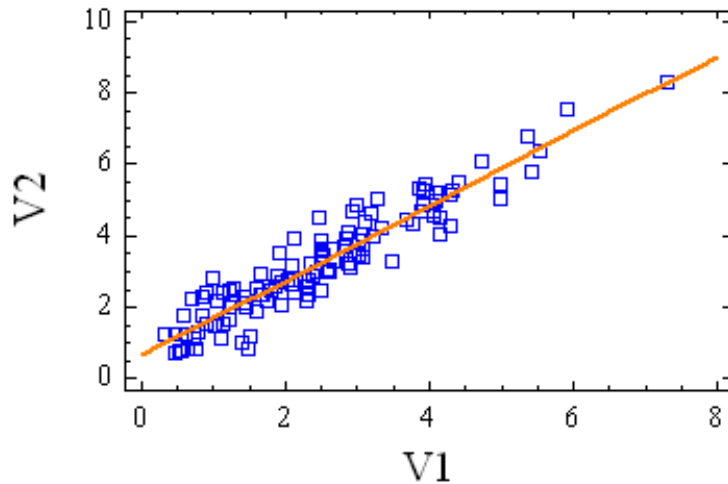


The absence of linearity is clear.  
In this case the regression line is not an useful  
information for prediction.

## Example

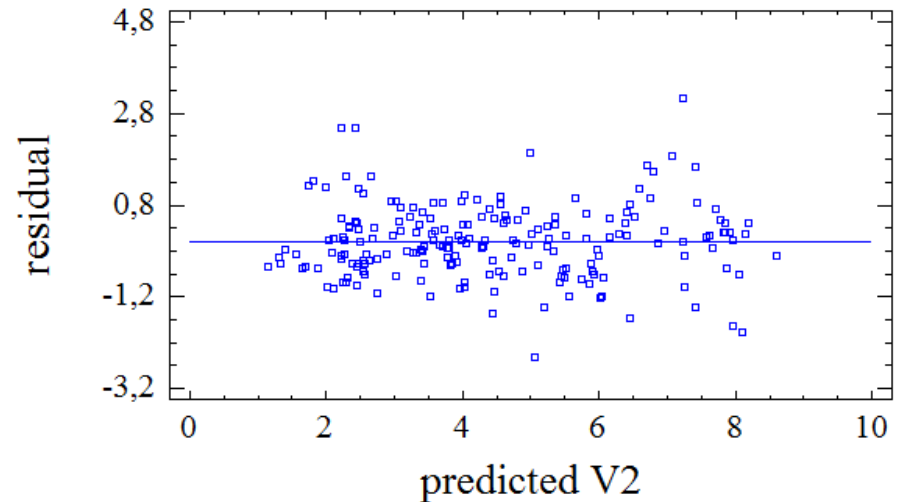
The variable  $V_1$  is the wind speed registered in location 1, while the variable  $V_2$  is the speed registered at the same time in location 2. There are a total of 115 pairs of measures.

Plot of Fitted Model



$$\hat{V}_2 = 0.657 + 1.045 \times V_1$$

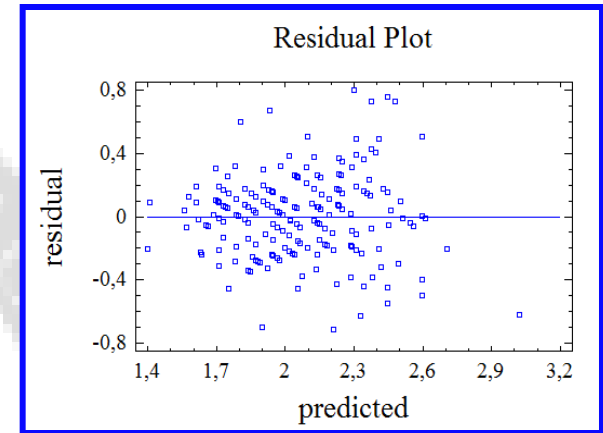
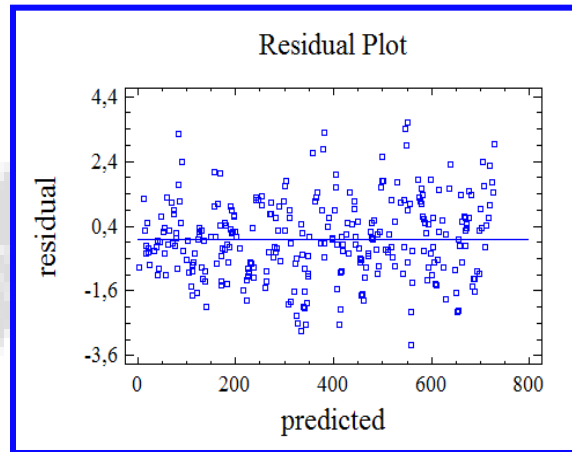
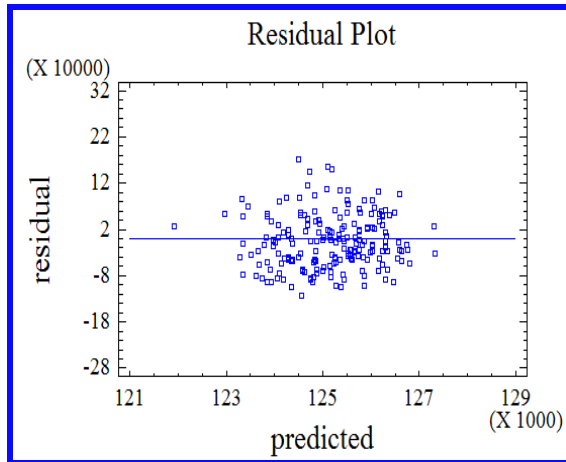
Residual Plot



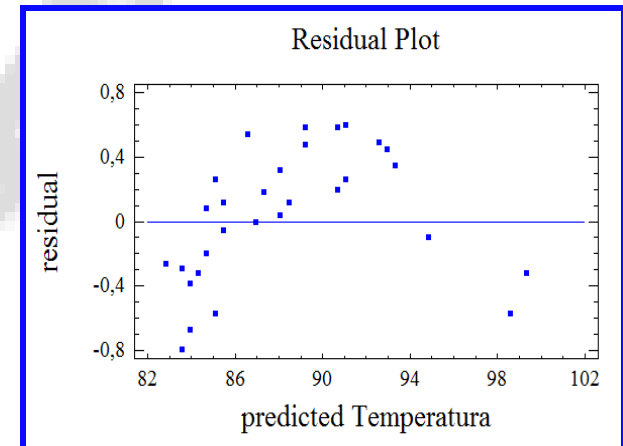
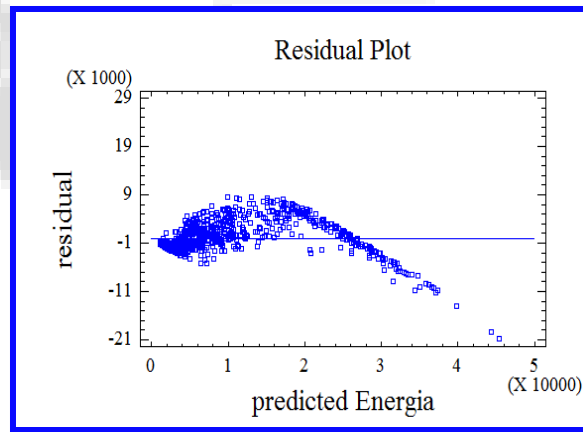
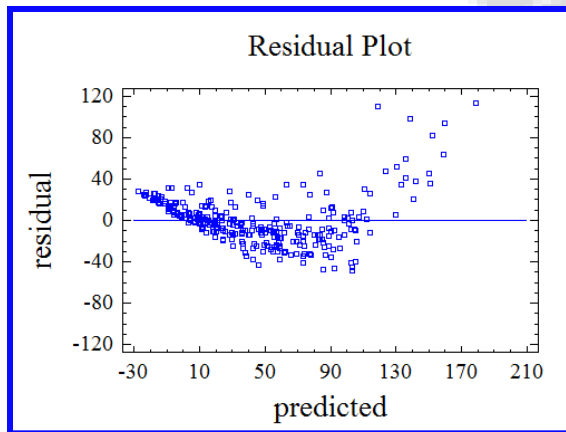
These residuals do not show any clear structure.

This means that the linear model is adequate

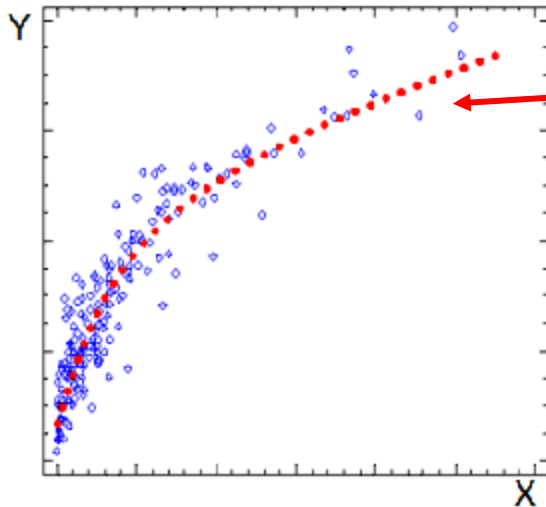
These residual graphs DO be acceptable



These residual graphs DO NOT be acceptable

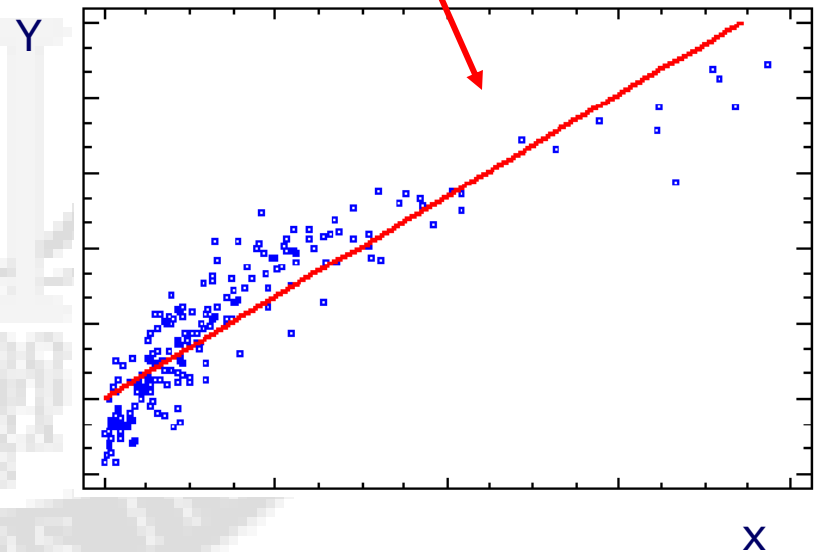


Some types of no linearity can be corrected by transforming the variable

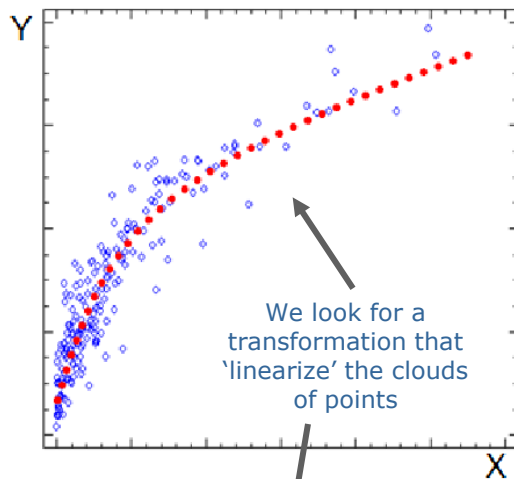


This curve is the one we would like to use to resume the relation...

... but the simple regression technique will give only this type of solutions

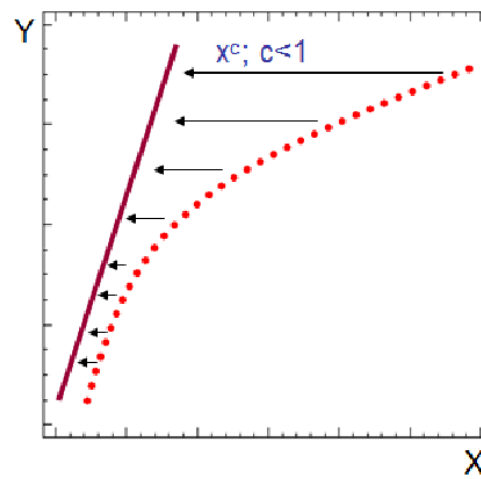


We look for an additional variables  
 $y^* = f(y)$  and  $x^* = g(x)$   
such that between them there is a  
linear relation



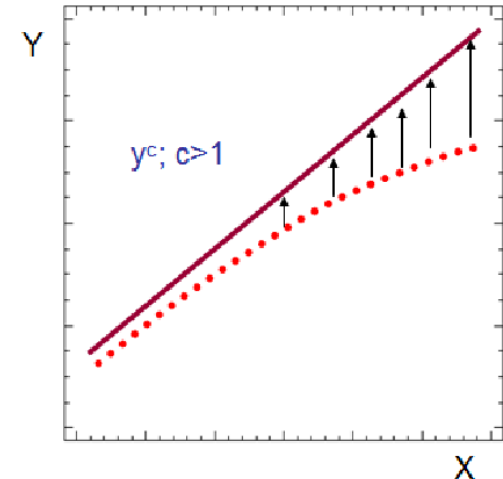
$$y = a + bx^c$$

$$y^c = a + bx$$



If  $c < 1$ , the biggest values are reduced the most. In this case applied to X we will 'linearize' the curve.

$$y = a + bx^c$$

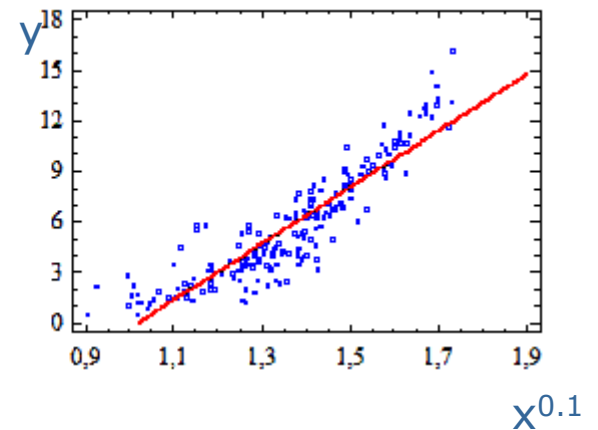
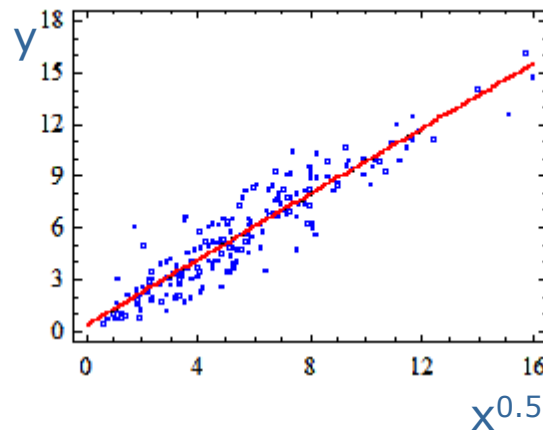
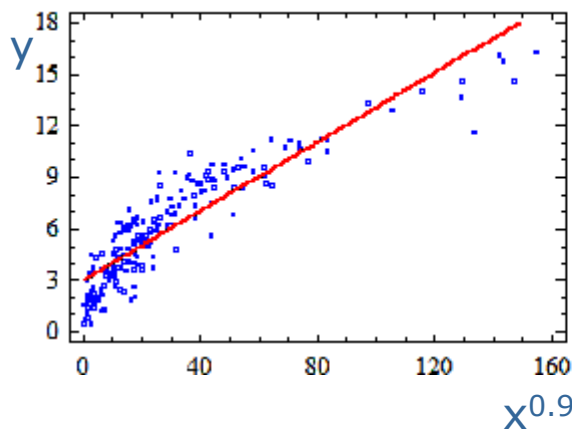


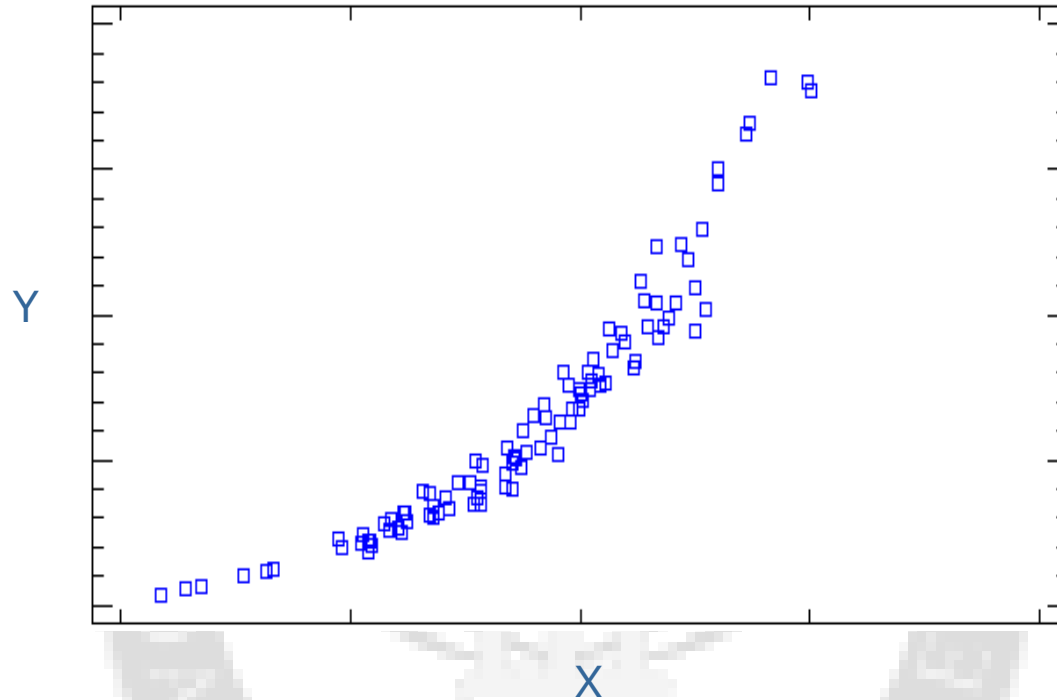
If  $c > 1$ , we get the opposite effect: the smallest values are expanded the most. In this case applied to Y we will 'linearize' the curve as well.

$c=0.9$  Insufficient

$c=0.5$  Perfect!

$c=0.1$  Too much





Assuming positive values. How the cloud of point will change if we use the following transformations:

- $y^2$
- $x^2$
- $y^{0.5}$
- $\log(y)$

## Evaluation of the regression

The regression to predict **y** starting from the value **x** will be good if:

1. The relation between x and y is linear ← Graphs
2. The linear relation is strong (big correlation)

Correlation and  
Square of correlation  
coefficient  $R^2$

$$R^2 = 1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

- Between 0 y 1
- $R^2 = \text{corr}(x, y)^2$
- The square of correlation ( or coefficient of determination) tells us which proportion of the dispersion of the dependent variable **y** is used to compute the regression line