

# Predictive Methods

*K. Gibert<sup>(1)</sup>*

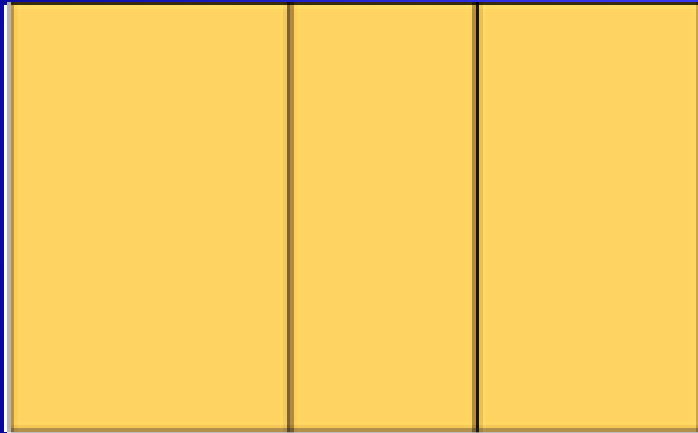
*<sup>(1)</sup>Department of Statistics and Operation Research*

*Knowledge Engineering and Machine Learning group  
Universitat Politècnica de Catalunya, Barcelona*

# Modelling

***Cognition***

Socio-econ. Opinions Products

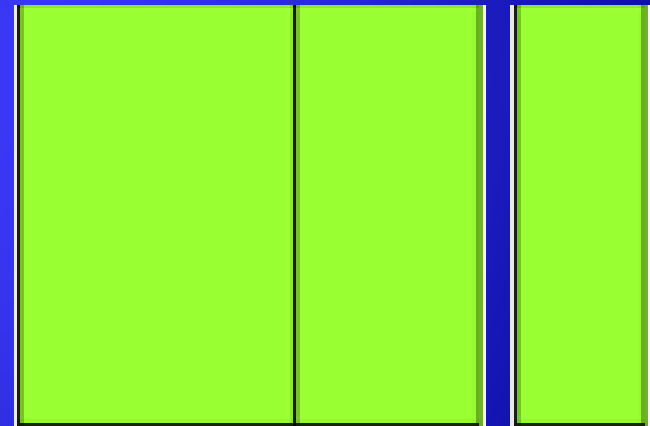


**Data to explore**

***Re-Cognition***

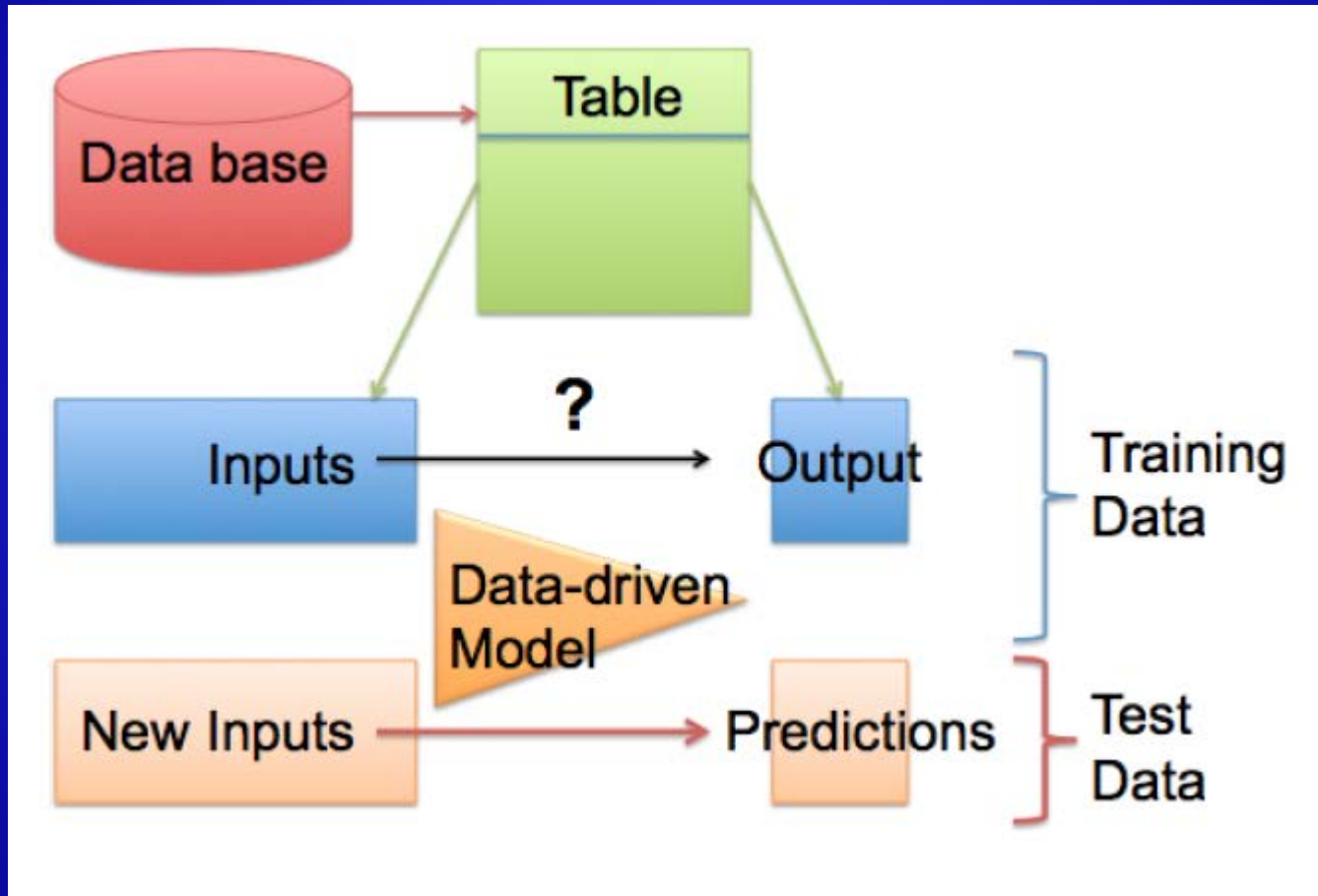
Inputs

Output(s)



**Data to modelize**

# Supervised Learning



# Supervised learning tasks

## DM goals [Fayyad et al., 1996]

- **Classification** – labeling a data item into one of several predefined classes (e.g. classify the type of credit client, “good” or “bad”, given the status of her/his bank account, credit purpose and amount);
- **Regression** – estimate a real-value (the *dependent variable*) from several (*independent*) attributes (e.g. predict the price of a house based on its number of rooms, age and other characteristics);

- **Classification:** Decision Tree, Random Forest, Classification Rules, Linear Discriminant Analysis, Naive Bayes, Logistic Regression, Neural Networks (MLP, RBF), SVM, ...
- **Regression:** Regression Tree, Random Forest, Multiple Regression, Neural Networks (MLP, RBF), SVM, ...

# Statistical Modelling

$$\text{Data} = \text{Fit} + \text{Error}$$

- Fit:
  - Structural
  - Law governing the phenomenon
  - Analytic Function
- Error:
  - Random
  - Variability around Fit (null expectation)
  - Probabilistic model

# Statistical models

- Determine the family of fits:

- Linear

- Quadratic

- Exponential

- .....

- Determine the law of error:

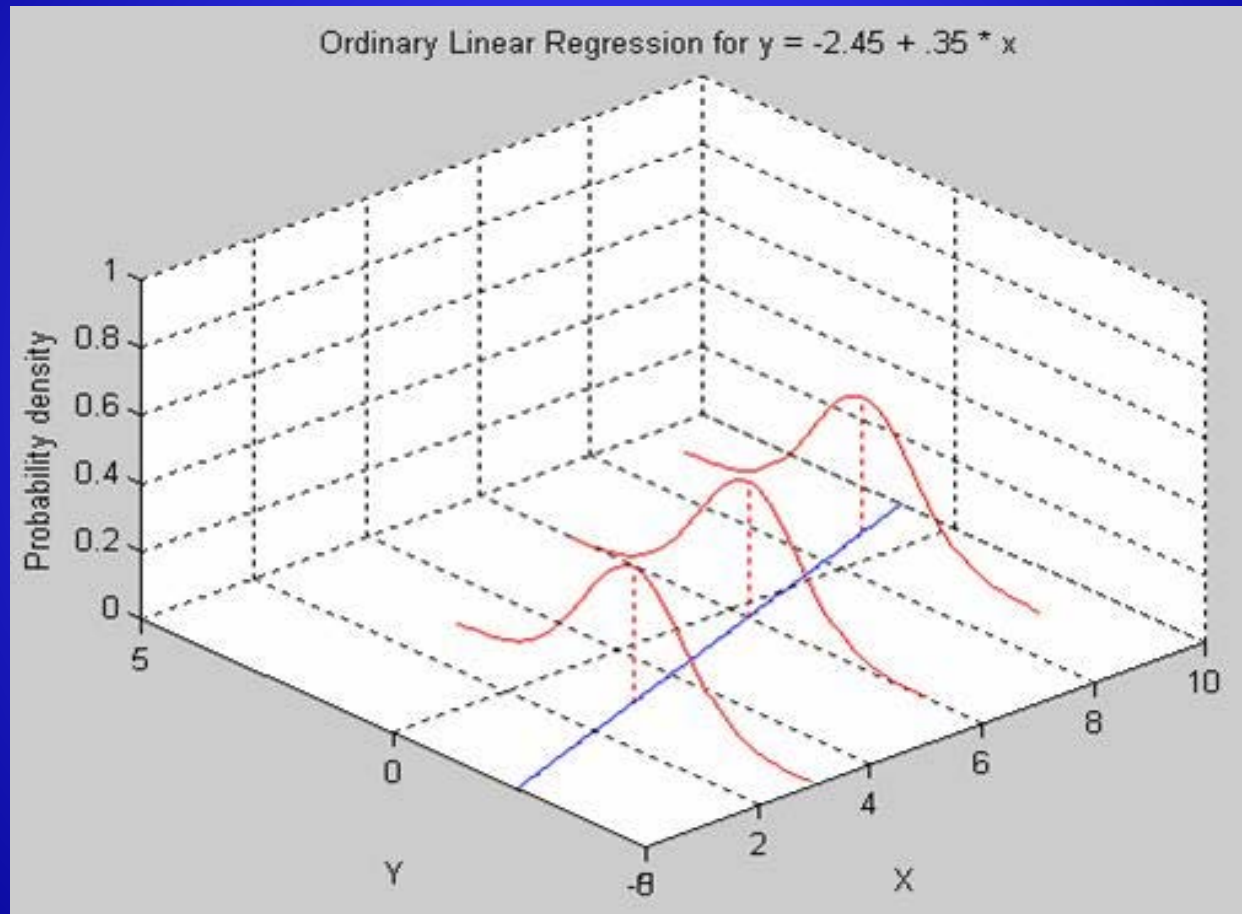
- Normal

- Poisson

- Binomial....



# E1 Normally Distributed Error



# Linear Multiple regression

- *Fit= linear; Error=Normal and centered*

- *Formalization:  $I=i:n$  observations*

*Y: Response variable*

*$X_1 \dots X_K$  : Explanatory Variables*

Find  $\beta_0 \dots \beta_K$  such that

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_K X_K + \varepsilon$$

- *Assumptions:*

- *Linearity:  $E(Y | X=x) = \mu_{y|x} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_K X_K$  ;  $E[\varepsilon] = 0$*

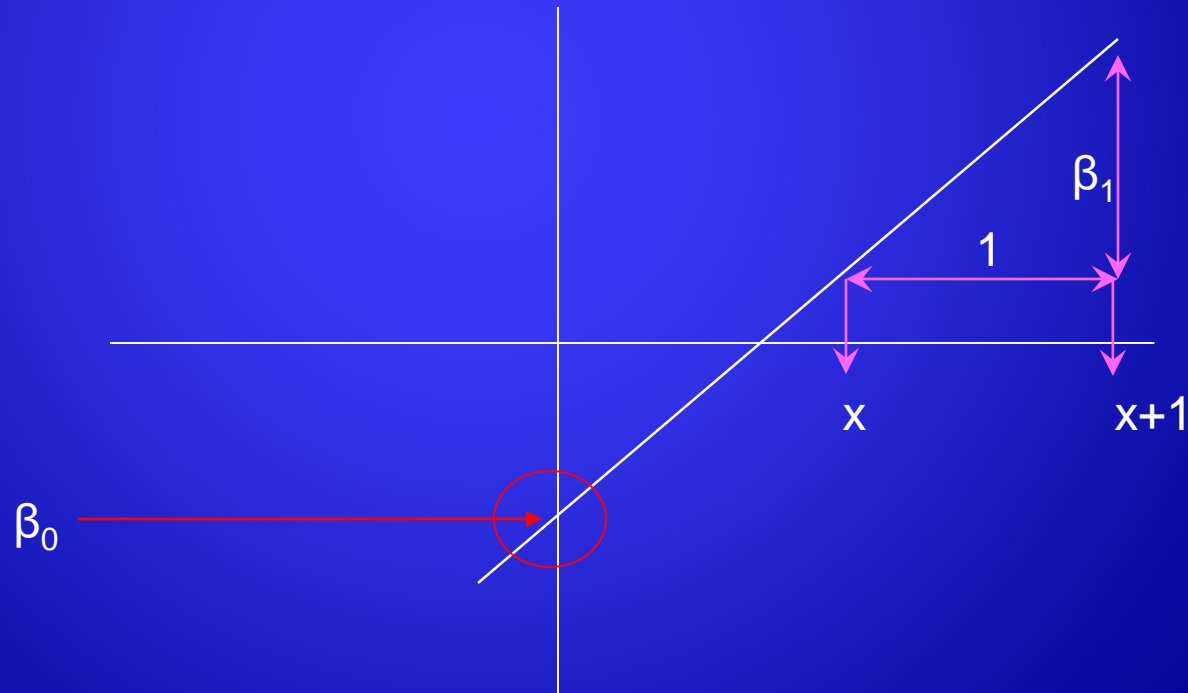
*Population regression line*

- *Normality:  $\varepsilon_1, \dots, \varepsilon_n \sim \mathcal{N}(0, \sigma_i)$ ,  $i=1:n$*
- *Homokedasticity:  $\text{Var}[\varepsilon_i] = \sigma^2$  for all  $i$*
- *Independence:  $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$  for all  $i, j$*

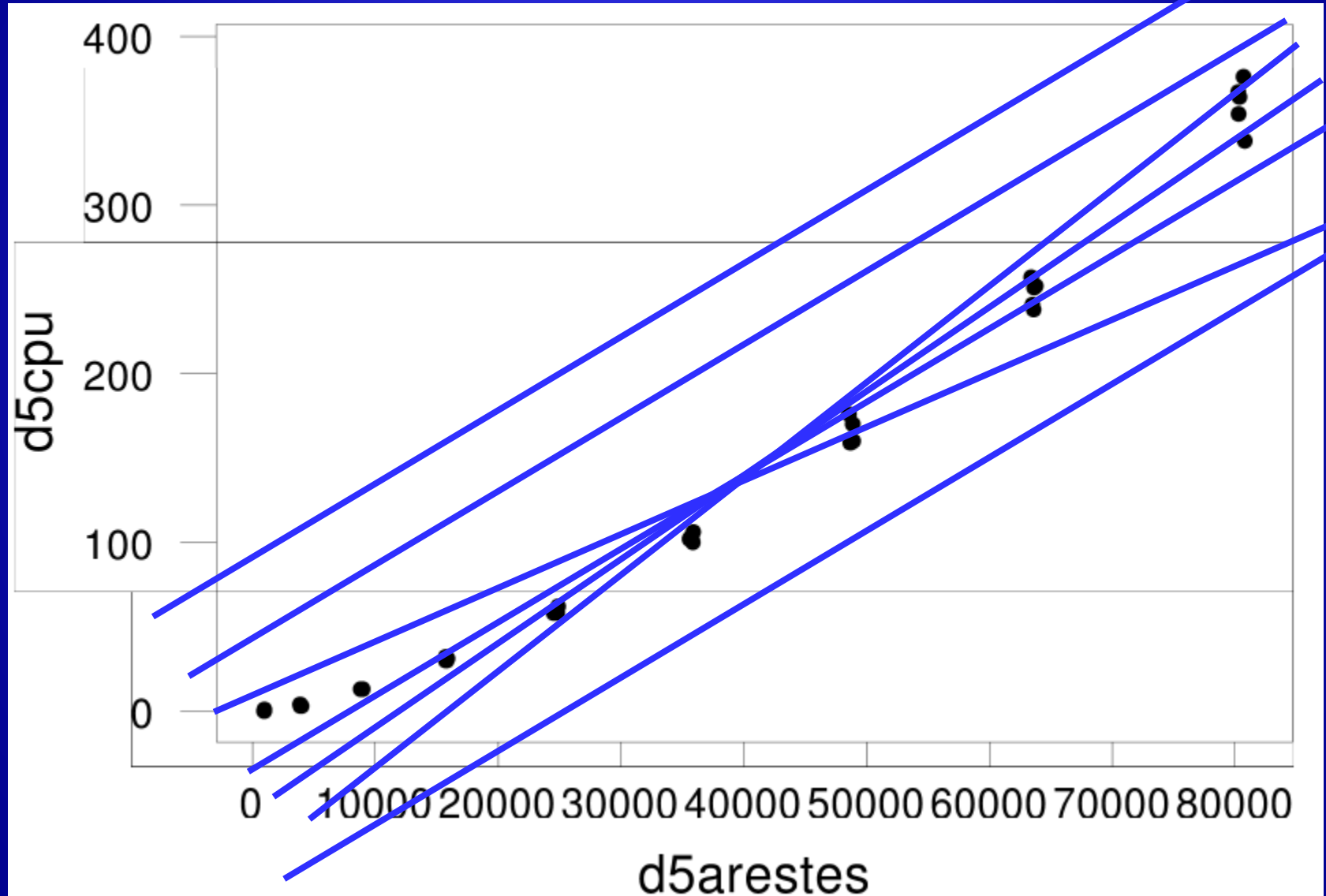


# What is “Linear”?

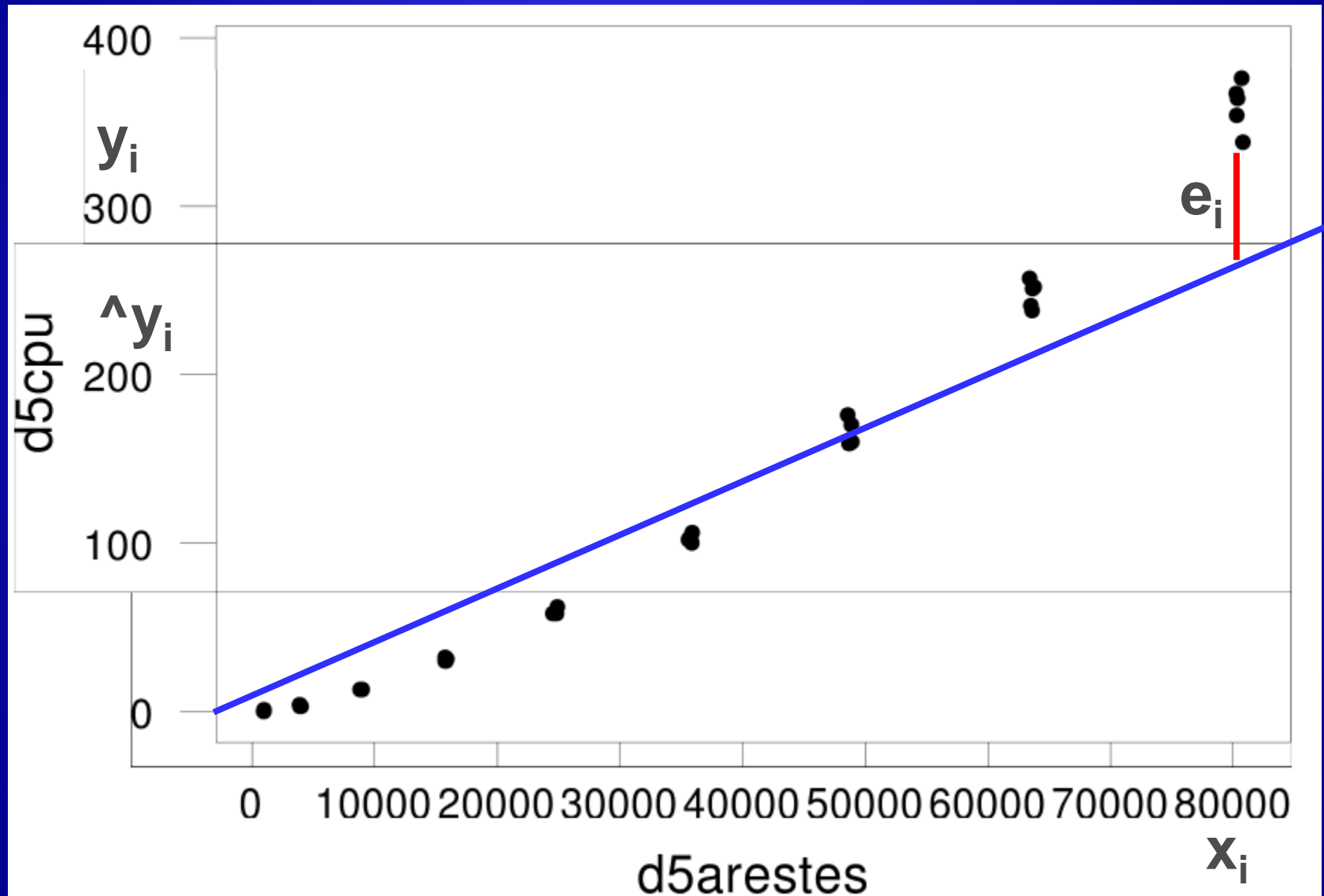
- Remember this:
- $Y = \beta_0 + \beta_1 X$  (simple linear regression, one  $X$ )



- Real case: Experimental CPU time of a graph treatment algorithm vs graph size



- Real case: Experimental CPU time of a graph treatment algorithm vs graph size



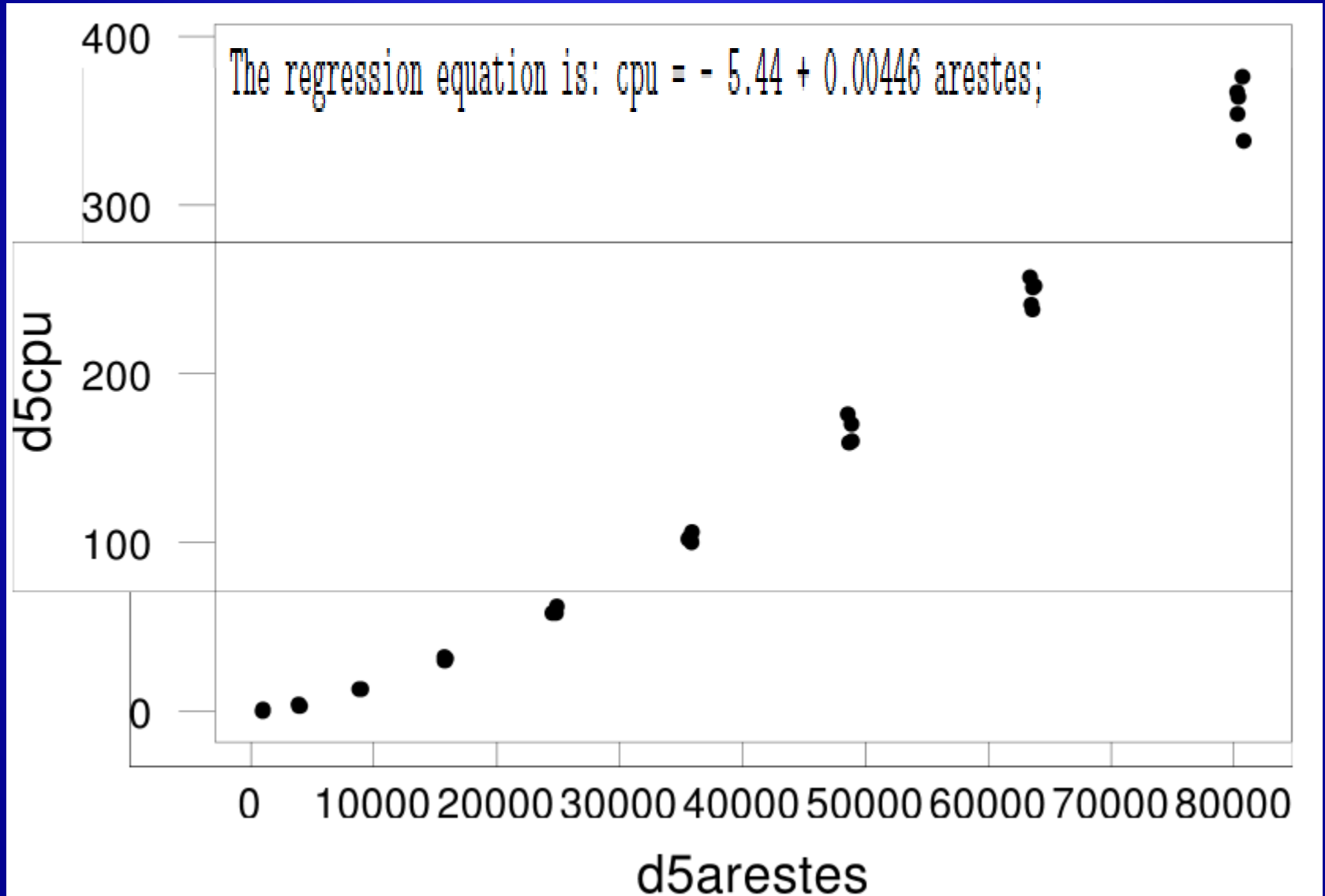
# Minimum Least Squares solution

Find  $\hat{\beta}_0, \hat{\beta}_1$  such that  $\min_{\hat{\beta}_0, \hat{\beta}_1} \sum_{\forall i} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$

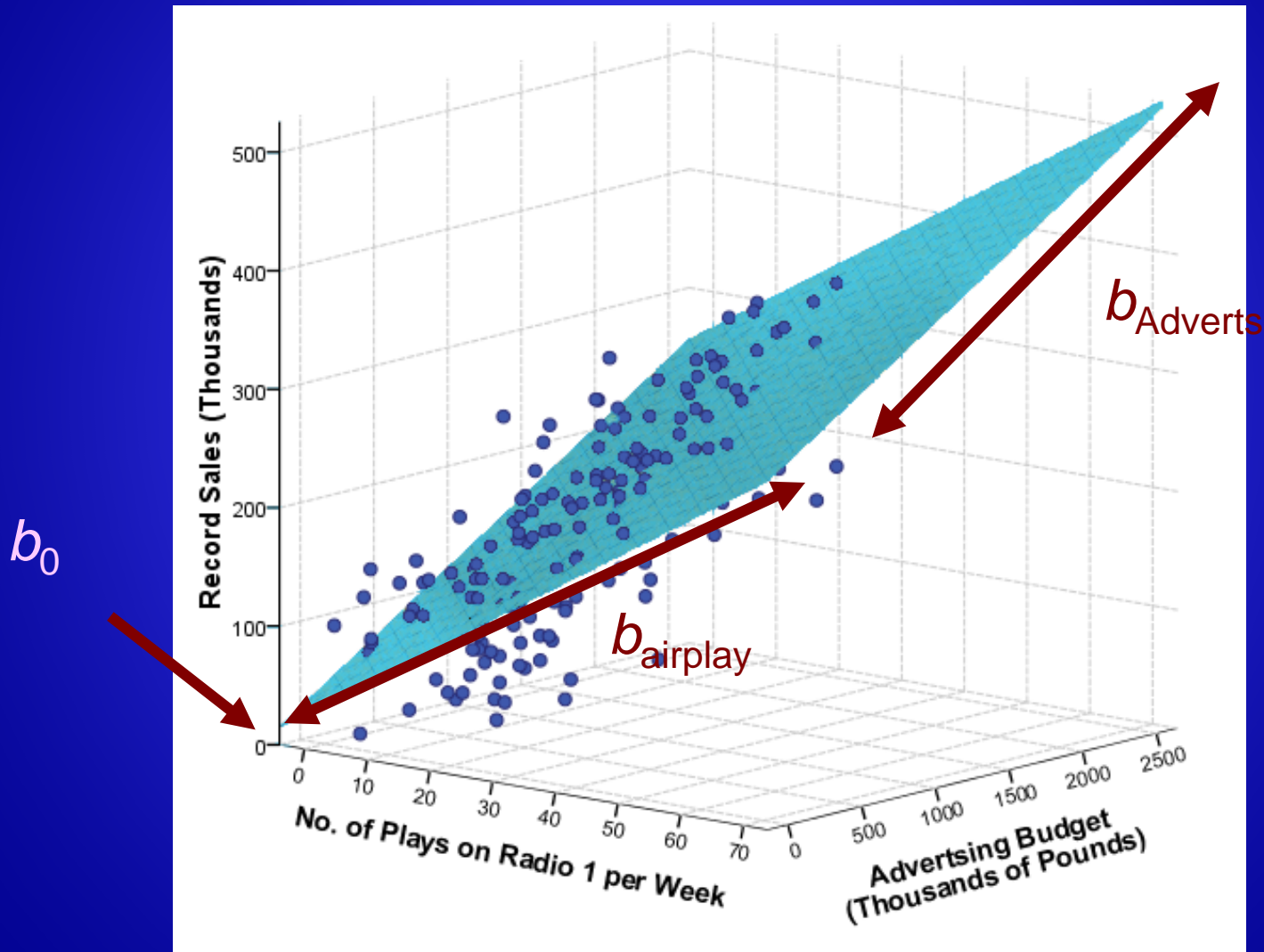
$$\hat{\beta}_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}$$

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

- Real case: Experimental CPU time of a graph treatment algorithm vs graph size



# The Model with Two Predictors





# Matricial formulation

$$\text{Regression fit criterion: } \min_r E \left[ \left( y_i - r(x_{i1}, \dots, x_{ip}) \right)^2 \right]$$
$$r(x_{i1}, \dots, x_{ip}) = E \left[ y_i \mid x_{i1}, \dots, x_{ip} \right]$$

$$E \left[ y_i \mid x_{i1}, \dots, x_{ip} \right] = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i$$

Estimation of coefficients

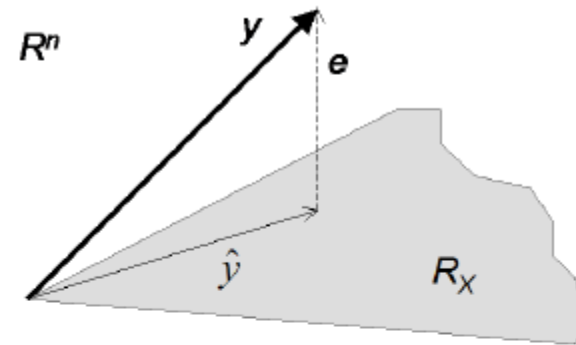
$$y_i = b_0 + b_1 x_{i1} + \dots + b_p x_{ip} + e_i$$

In matrix notation

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_p \end{bmatrix} + \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix} \equiv y = Xb + e = \hat{y} + e$$

# Geometric interpretation

$$y_i = \hat{y}_i + e$$
$$\begin{matrix} y & \hat{y} & e \\ \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} & = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} + \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix} \end{matrix}$$



$$\hat{y}_i = b_0 + b_1 x_{i1} + \dots + b_p x_{ip}$$

$$\text{Criterion: } \min_{b_0, \dots, b_p} \sum_{i=1}^n (e_i)^2 = \|e\|^2$$

$$\langle \hat{y}, e \rangle = \langle \hat{y}, y - \hat{y} \rangle = 0$$

$$\hat{y} = Xb, \quad b'X'y - b'X'Xb = 0$$

$$b = (X'X)^{-1} X'y$$

# Validation

- Technical Assumptions
  - normality, linearity, independence, homokedasticity
  - Tools
    - Graphical residuals analysis
    - Influence-point indicators (hi)
- Quality:
  - $R^2$  (determination coefficient): goodness, reliability
  - $s^2$ : noise, precision
  - Both guarantee generalizability (only interpolation)

# Quantify Goodnes of model

$$s^2 = \hat{\sigma}^2 = \frac{\sum_{i=1}^n (e_i)^2}{n-2}$$

Estimates the variance of residuals

The biggest, the worst the model, more impresice predictions



**Non-standardized**

# Quantify Goodnes of fit

$R^2$  : proportion of explained variance

$SStotal = V(Y)$  variance of response variable

Decomposition:  $SStotal = SS_{explainedByModel} + SS_{error}$

Dividing all sides by  $SStotal$ :

$$R^2 = \frac{SS_{explainedByModel}}{SStotal} = 1 - \frac{SS_{error}}{SStotal}$$

# Quantify Goodnes of model

$$SSTotal = V(Y) = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}$$

$$SSExplainedByModel = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{k-1}$$

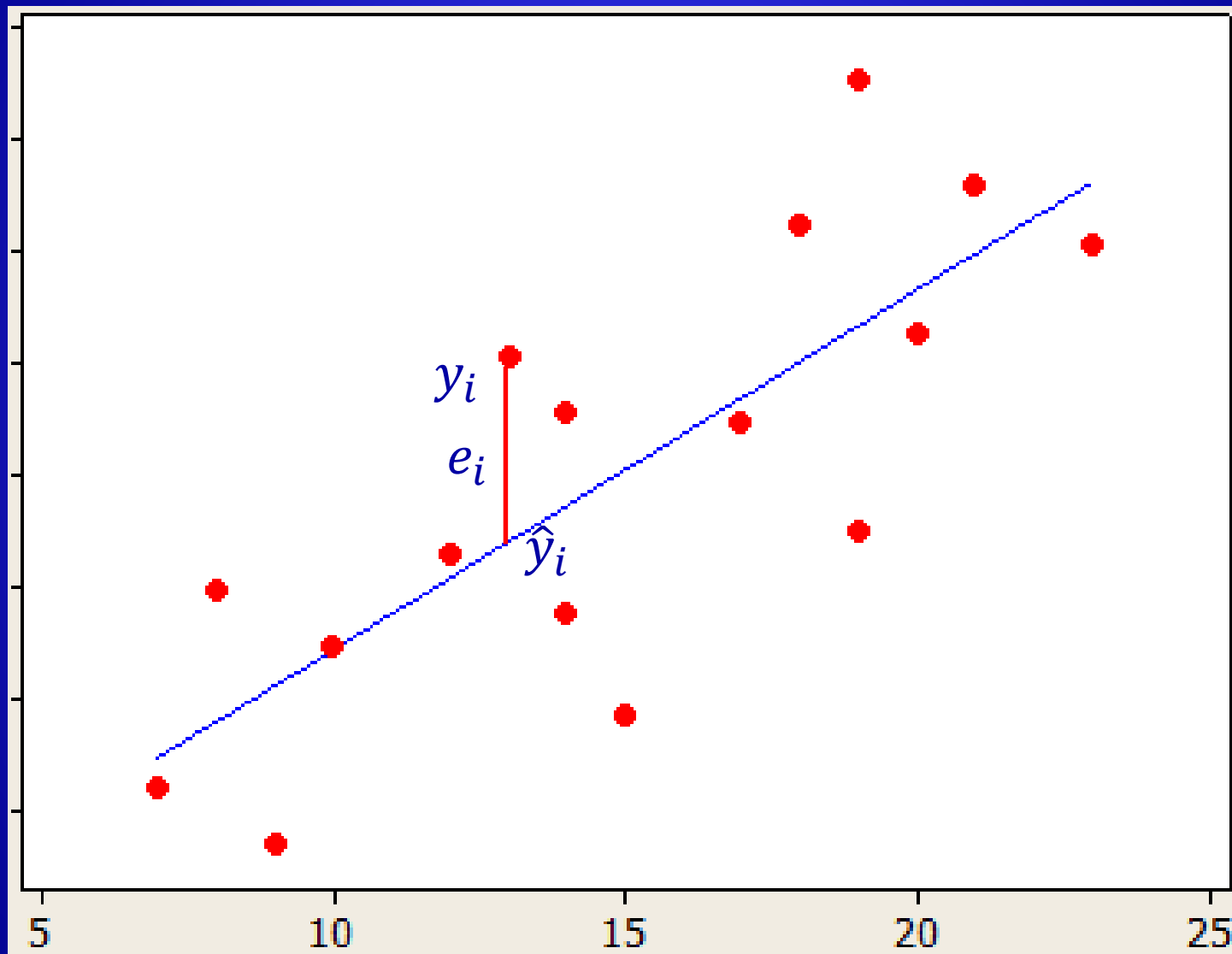
$$SSError = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-k}$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \dots + \hat{\beta}_k x_{ki}$$

$$\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$$



# The residuals



# Quantify Goodnes of model

$$SSTotal= V(Y) = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}$$

$$SSExplainedByModel= \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{k-1}$$

$$SSError= \frac{\sum_{i=1}^n (e_i)^2}{n-k}$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \dots + \hat{\beta}_k x_{ki}$$

$$\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$$

# Quantify Goodnes of model

$R^2$  = proportion of explained variance

$$R^2 = 1 - \frac{SSError}{SSTotal}$$

$$0 < R^2 < 1$$

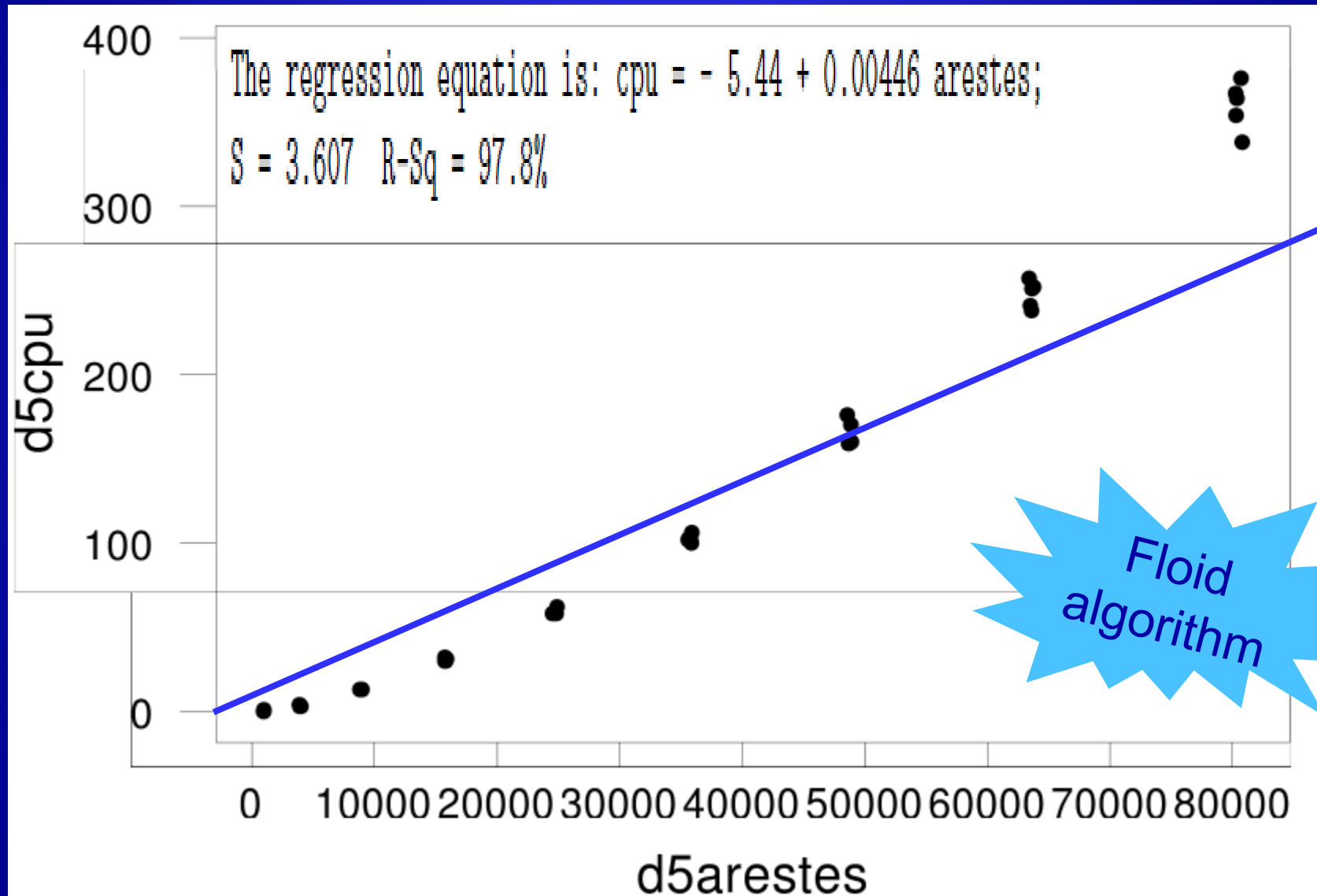
The biggest  $R^2$ , the better the model explains Y

For simple linear regression  $R^2 = \text{Corr}(Y, X)^2$



assume  
linearity

- Real case: Experimental CPU time of a graph treatment algorithm vs graph size



# Model inference

To test significance of the model

$$F = \frac{SSE_{ExplainedByModel}}{SSE_{Error}} \sim F_{(k-1, n-k)}$$

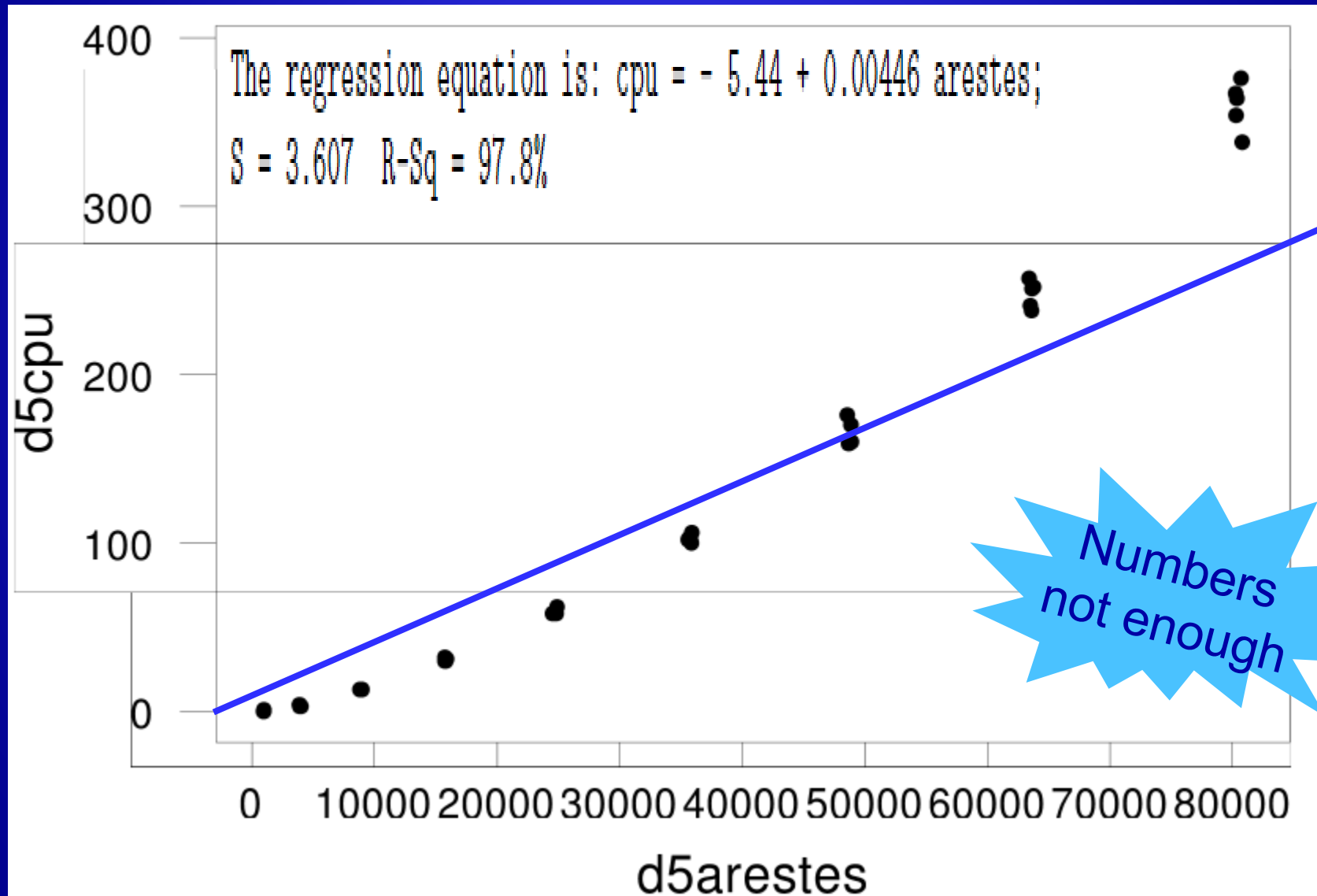
To test significance of a model term  $\hat{\beta}_k$

$$t_k = \frac{\hat{\beta}_k}{S_{\hat{\beta}_k}} \sim t_{n-K}$$

To test significance of a model term

**Both  
assume  
normality**

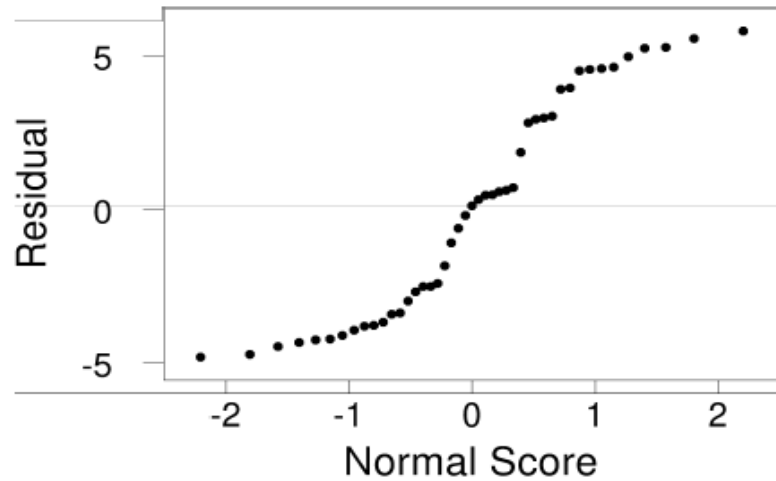
- Real case: Experimental CPU time of a graph treatment algorithm vs graph size



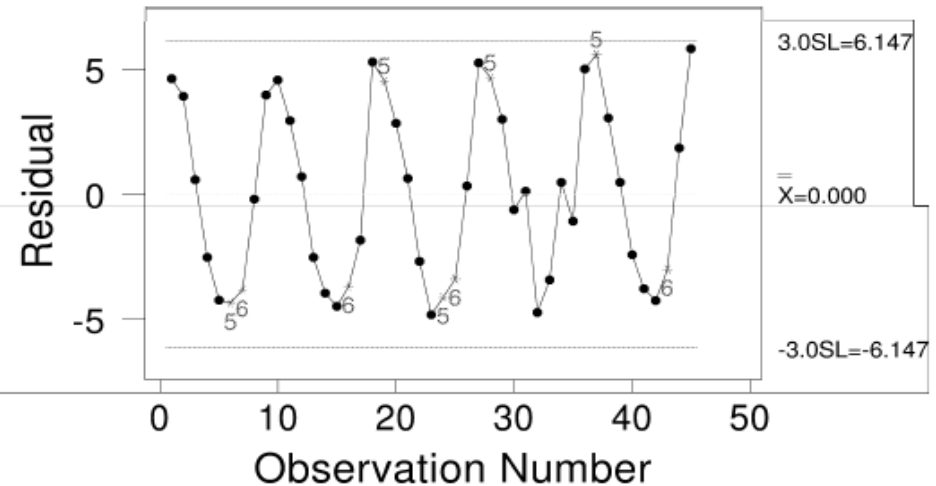


# Graphical residuals analysis

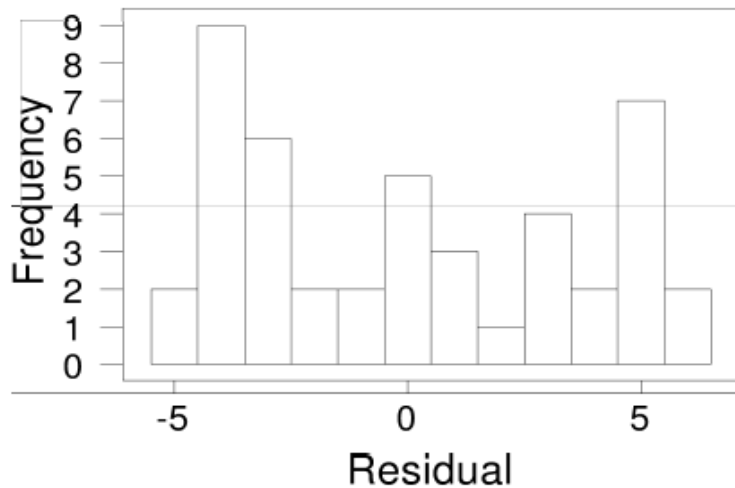
Normal Plot of Residuals



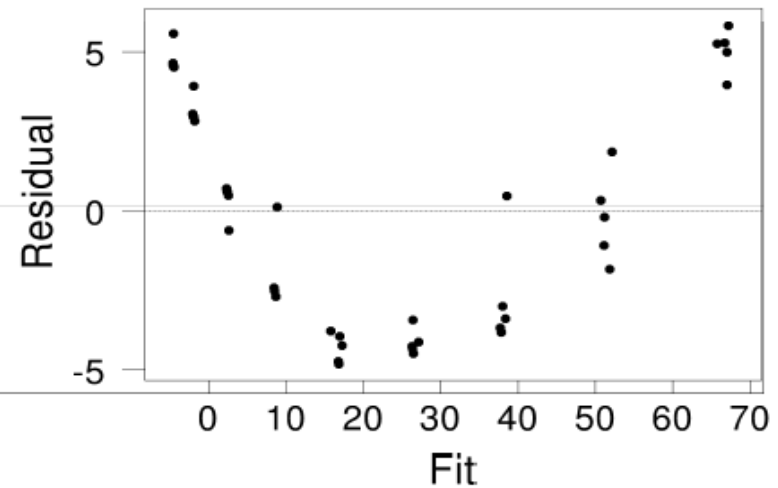
I Chart of Residuals



Histogram of Residuals



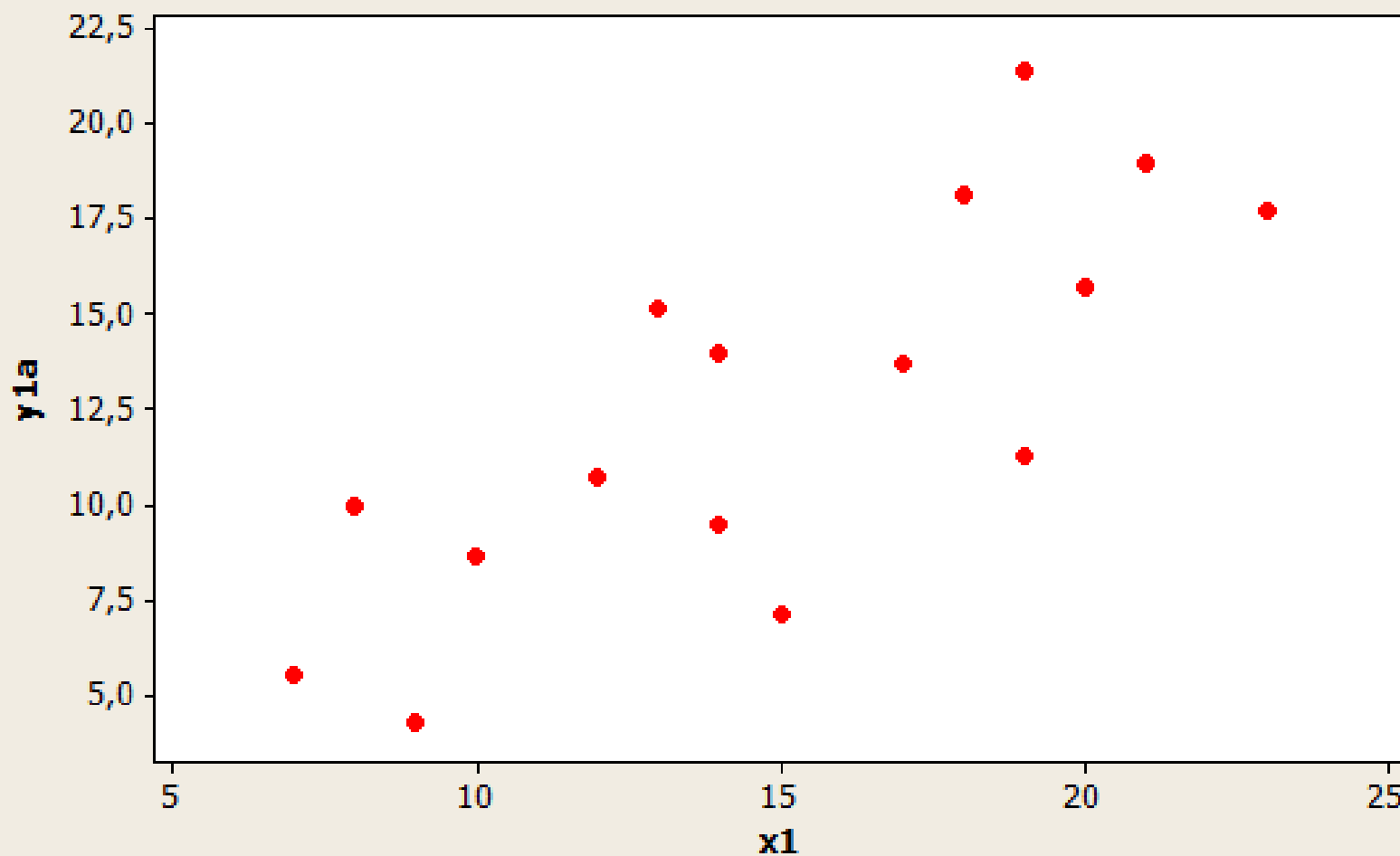
Residuals vs. Fits



# Regression

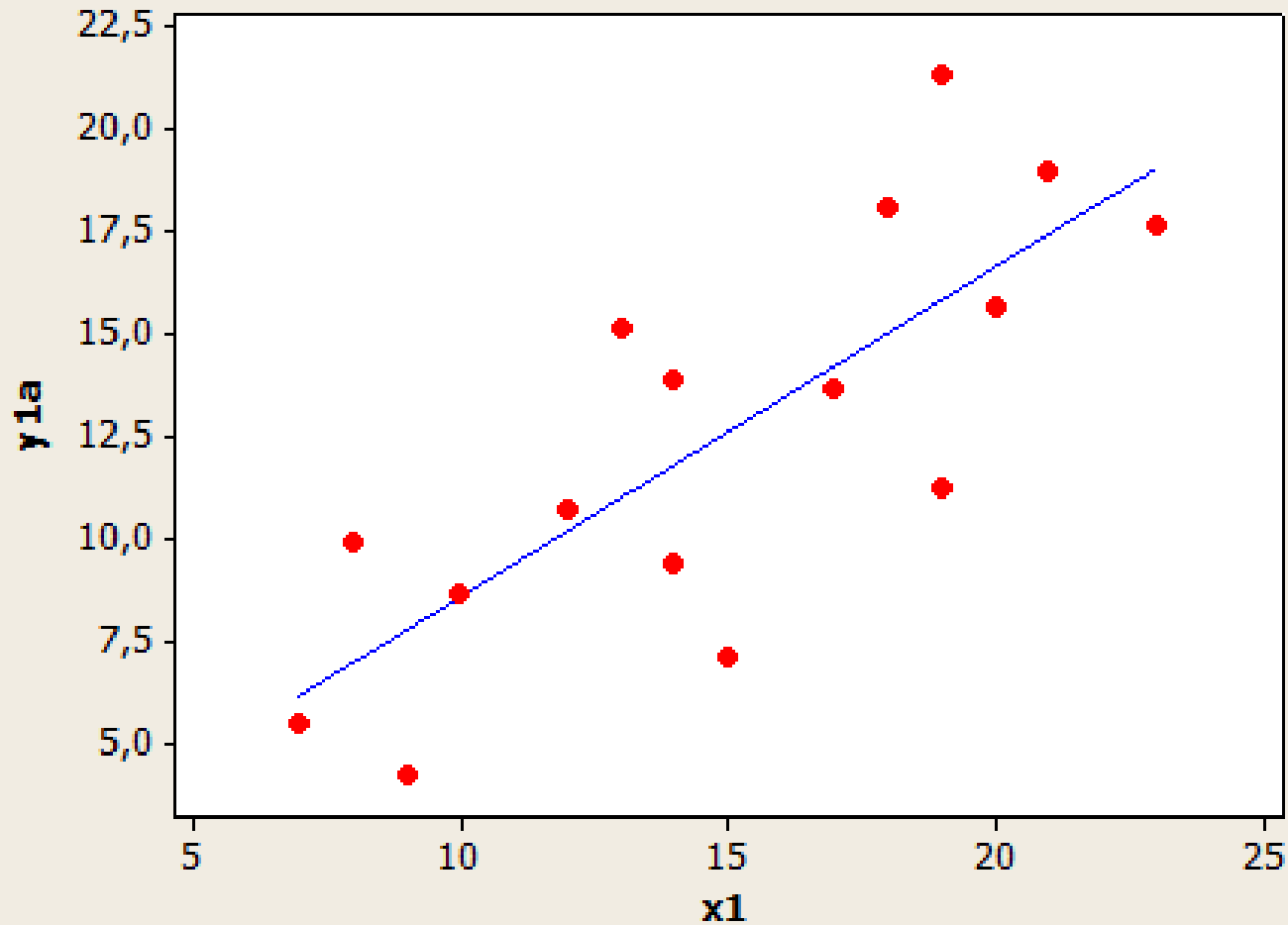
[Tomassone 56]

Scatterplot of y1a vs x1



### Fitted Line Plot

$$y1a = 0,522 + 0,8085 x1$$

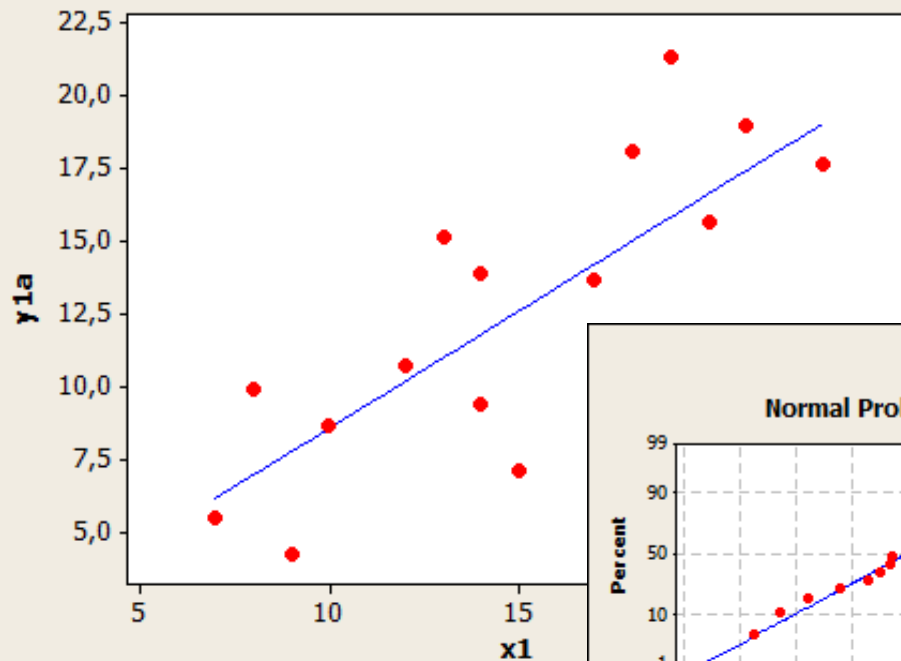


S	3,22314
R-Sq	61,7%
R-Sq(adj)	59,0%

# Regression

**Fitted Line Plot**

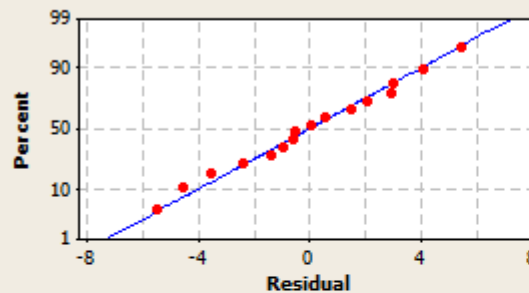
$$y1a = 0,522 + 0,8085 x1$$



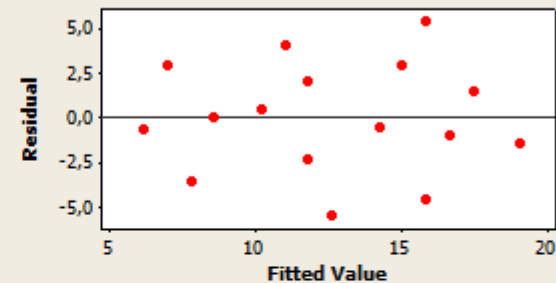
S	3,22314
R-Sq	61,7%
R-Sq(adj)	59,0%

**Residual Plots for y1a**

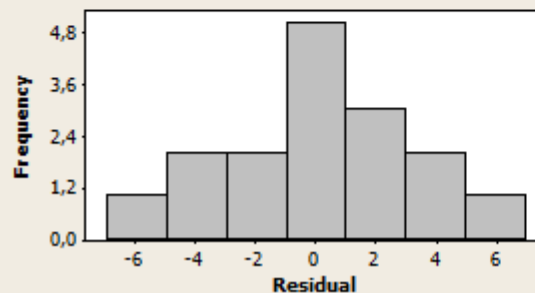
**Normal Probability Plot**



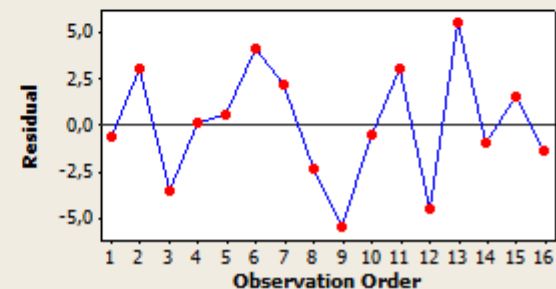
**Versus Fits**



**Histogram**

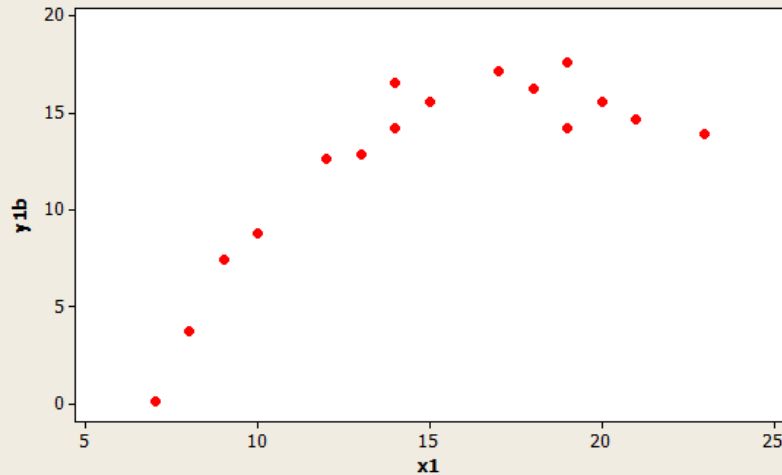


**Versus Order**

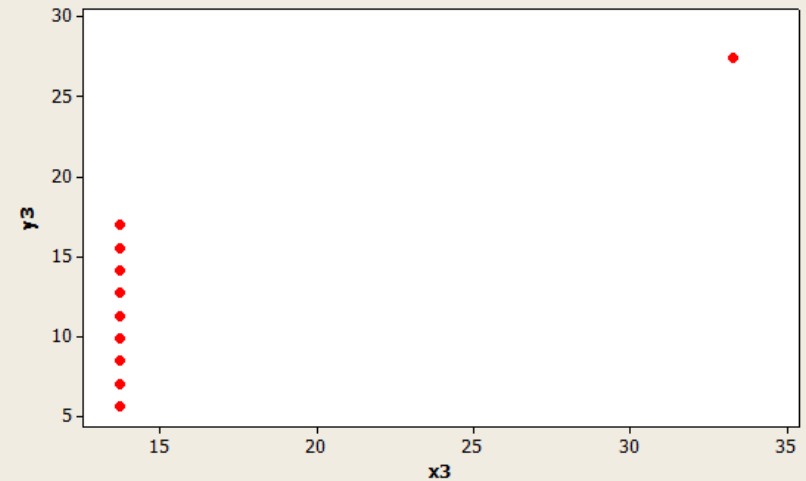


# Regression

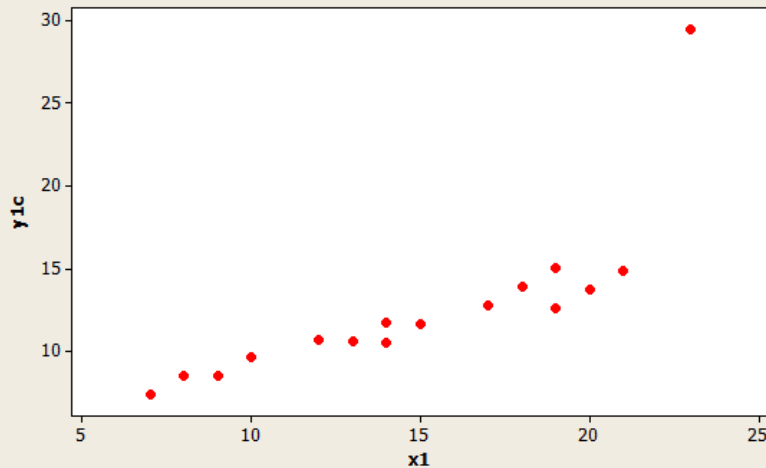
Scatterplot of y1b vs x1



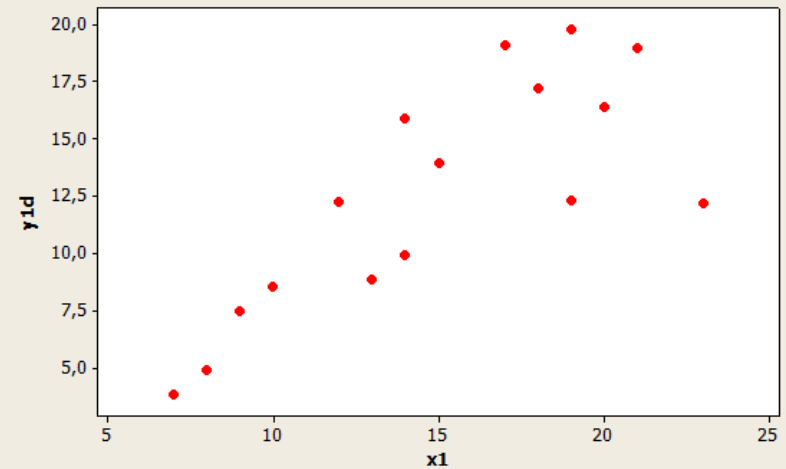
Scatterplot of y3 vs x3



Scatterplot of y1c vs x1

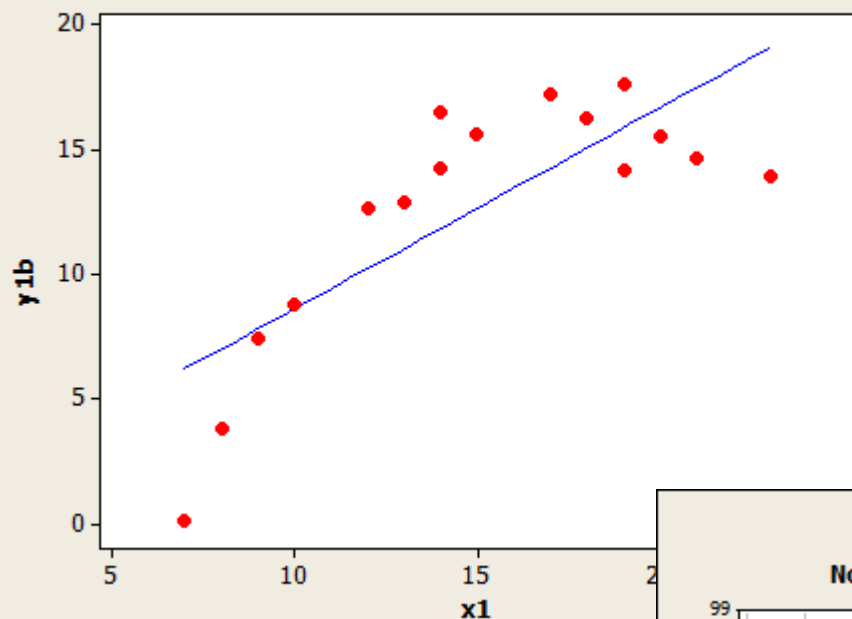


Scatterplot of y1d vs x1



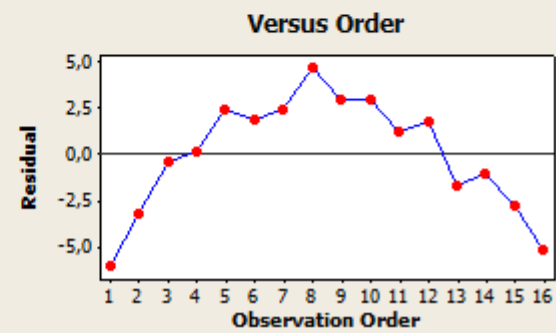
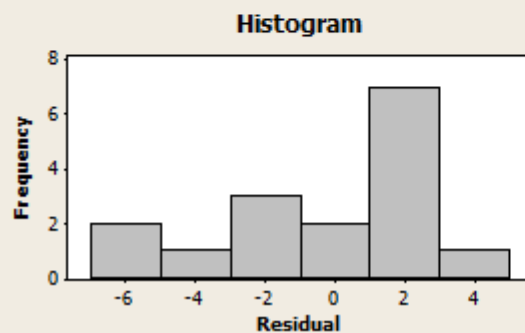
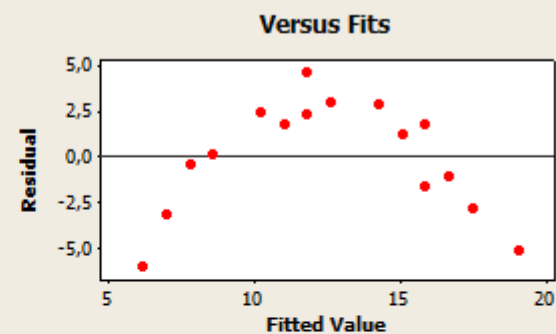
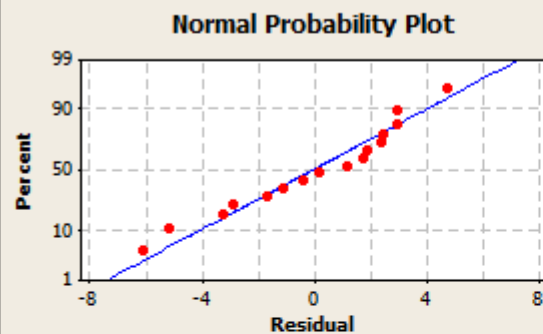
# Fitted Line Plot

$$y1b = 0,524 + 0,8085 x1$$



S	3,22655
R-Sq	61,7%
R-Sq(adj)	58,9%

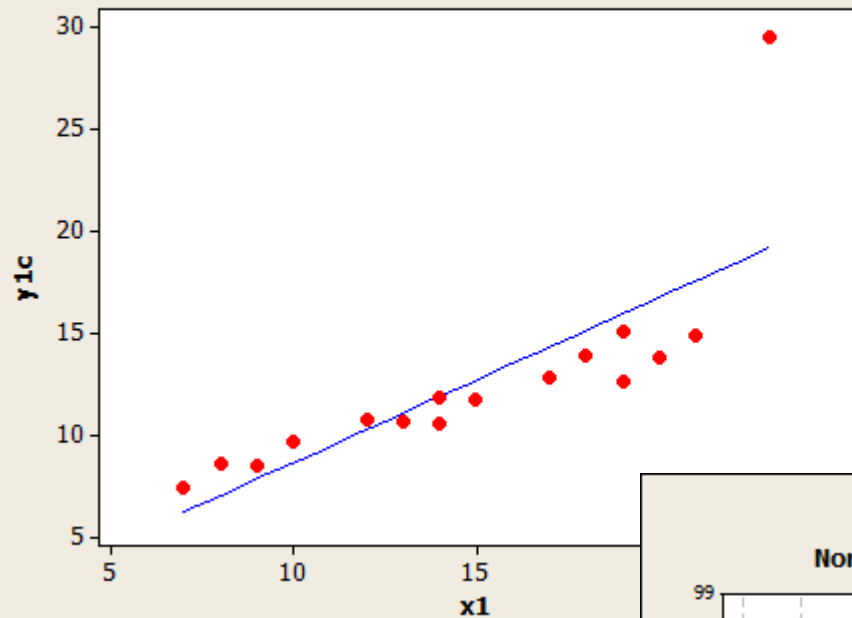
## Residual Plots for y1b





# Fitted Line Plot

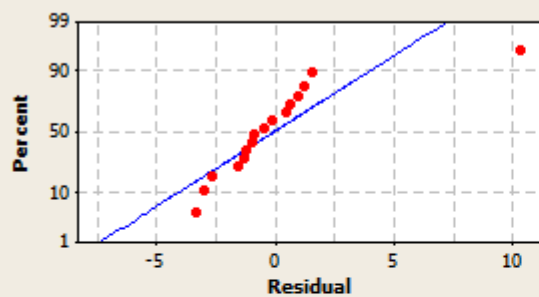
$$y1c = 0,520 + 0,8087 x1$$



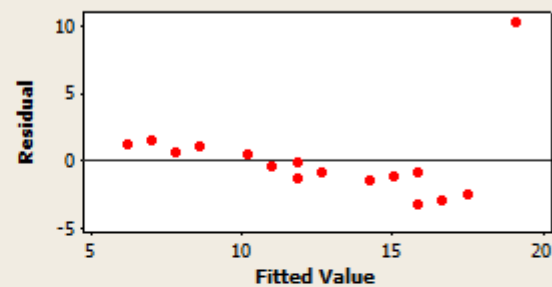
S	3,22553
R-Sq	61,7%
R-Sq(adj)	59,0%

## Residual Plots for y1c

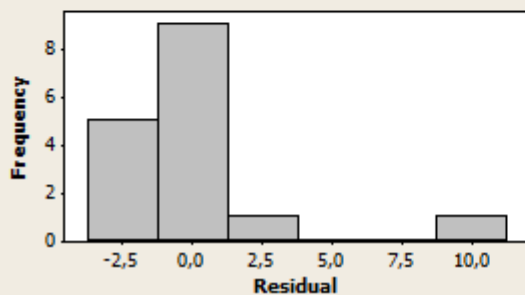
Normal Probability Plot



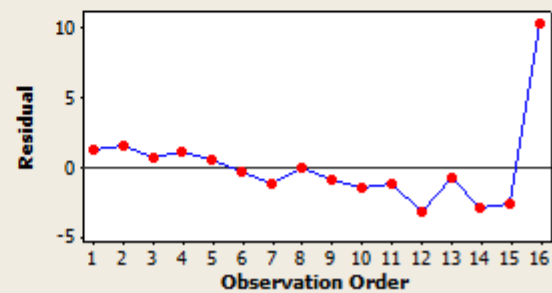
Versus Fits



Histogram

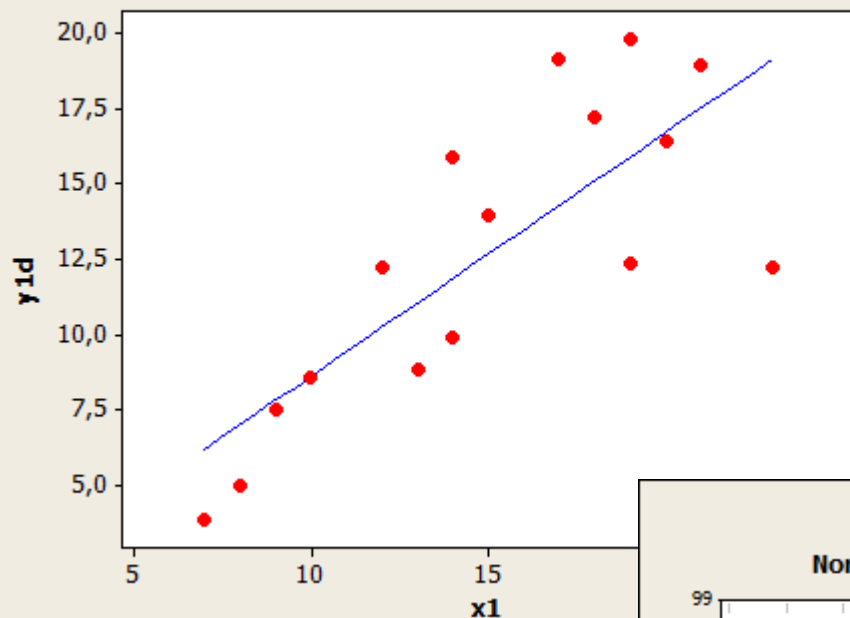


Versus Order



# Fitted Line Plot

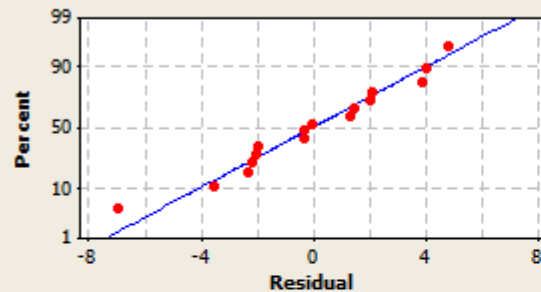
$$y1d = 0,520 + 0,8087 x1$$



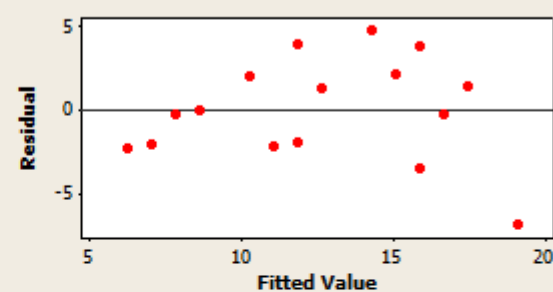
S	3,22559
R-Sq	61,7%
R-Sq(adj)	59,0%

## Residual Plots for y1d

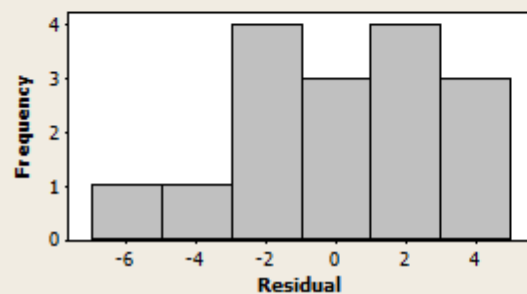
### Normal Probability Plot



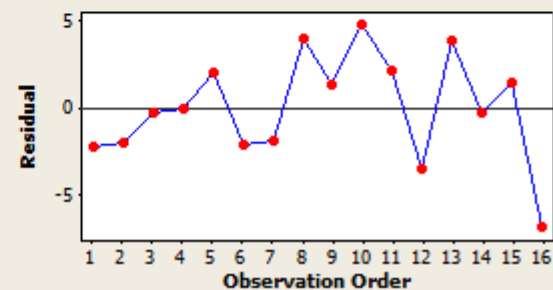
### Versus Fits



### Histogram

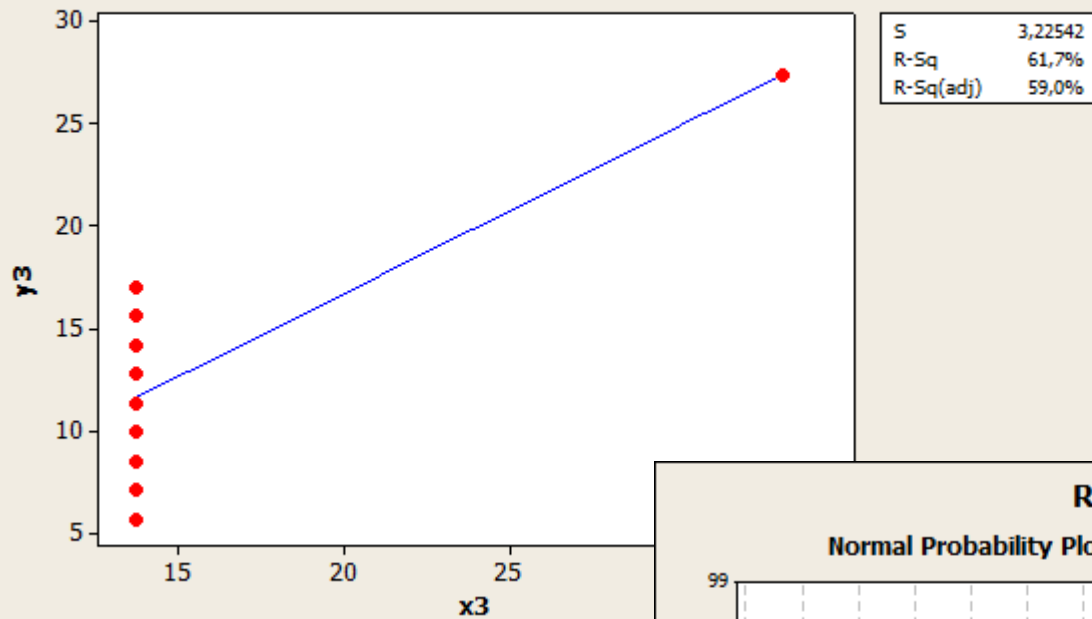


### Versus Order

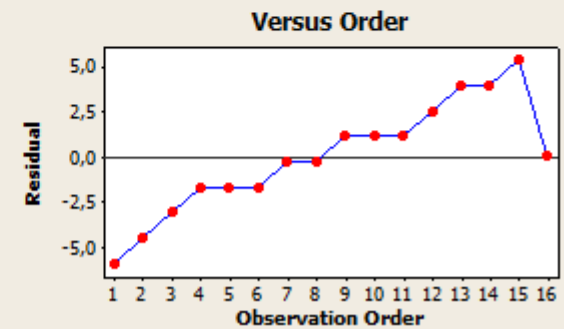
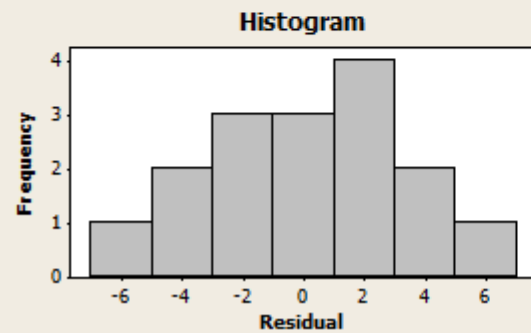
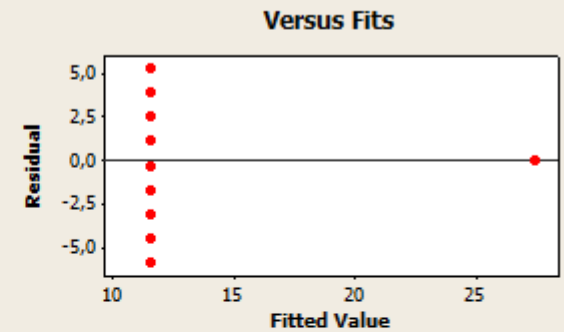
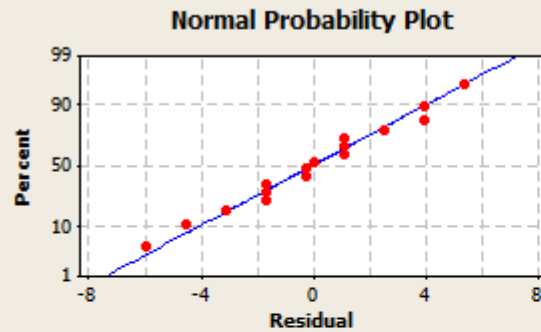


### Fitted Line Plot

$$y3 = 0,519 + 0,8087 x3$$



### Residual Plots for y3



# Going further

## Multiple linear regression

- Uses the same principles as simple linear regression
- Contains several regressions in the right hand side of regression equation

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_K X_K + \varepsilon$$

- All properties of simple linear regression hold, except one
  - Determination coefficient ( $R^2$ ) still measures goodness of fit
  - But  $R^2$  is not equal to the squared correlation coefficient in multiple regression
- Analysis of the residuals is done by every regressor  $X_k$

# Modelling strategy

- 1) plots between  $Y$  and all  $X_k$  to check the existence of linear association or the need of entering some  $X_k^2$  or some  $1/X_k$  (provided that they make sense)
- 2) Check p-values of all coefficients in the final regression equation. Repeat the model with all those with non significant p-value. Repeat till all coefficients look significant
- 3) With the model with all significant coefficients, write the equation of the final regression model
- 4) Make the graphical analysis of the residuals and verify all checks are passed
- 5) If some check fails,..... Proceed accordingly and repeat the model
- 5) Report final equation and  $R^2$  . Also include  $R^2$  value of the very first model
- 6) report all intermediate work done to move from the very first model to the very last and all decisions made in the meanwhile

# Going much further

- ANCOVA: to introduce qualitative variables
- Interaction terms to introduce multiplicative models
- Polynomic regression to estimate higher order polynomial functions
- General Linear Model (common formulation for simple/multiple linear regression, ANOVA and ANCOVA)
- Generalized linear models: common formulation for an extension of families of models:
  - Linear,
  - Poisson
  - Logit.....
- Non linear relationships: LOESS (Locally Weighted Least Squares Regression), uses more local data to estimate the model. It uses a 'nearest neighbors' method to smooth data.
- Complex functions: Artificial Neural Networks



# Supporting materials

Basic linear regression:

<https://www.scribd.com/doc/304772056/BI-147-Web> Chapter 2

[https://www.google.es/url?sa=t&rct=j&q=&esrc=s&source=web&cd=9&ved=2ahUKEwiQ7u\\_nx\\_noAhWE3oUKHYIQCo8QFjAlegQIChAB&url=https%3A%2F%2Fwww.parisnanterre.fr%2Fmedias%2Ffichier%2Fcours\\_regression\(2\)\\_1273084206969.pdf%3FID\\_FICHE%3D204222%26INLINE%3DFALSE&usg=AOvVaw3bXGbamKKITIn5ewLD33DS](https://www.google.es/url?sa=t&rct=j&q=&esrc=s&source=web&cd=9&ved=2ahUKEwiQ7u_nx_noAhWE3oUKHYIQCo8QFjAlegQIChAB&url=https%3A%2F%2Fwww.parisnanterre.fr%2Fmedias%2Ffichier%2Fcours_regression(2)_1273084206969.pdf%3FID_FICHE%3D204222%26INLINE%3DFALSE&usg=AOvVaw3bXGbamKKITIn5ewLD33DS)

Linear Regression with R:

<http://r-statistics.co/Linear-Regression.html>

Multiple linear regression (matricial notation)

<http://dept.stat.lsa.umich.edu/~kshedden/Courses/Stat401/Notes/401-multreg.pdf>

<https://www.google.es/url?sa=t&rct=j&q=&esrc=s&source=web&cd=2&ved=2ahUKEwji5-qkwvnoAhV0kFwKHRbiAscQFjABegQIAxAB&url=https%3A%2F%2Fwww.uv.es%2Furiel%2FChapter%25203%2520Slides.pdf&usg=AOvVaw0qzppaA8Aj6vRa79hrIBxr>

Moore, McCabe. "Craig (2012) Introduction to the Practice of Statistics." Chap 10

For the graphical analysis of the residuals:

Tomassone, Richard, Elisabeth Lesquoy, and Claude Millier. "La régression. Nouveaux regards sur une ancienne méthode statistique." (1983).