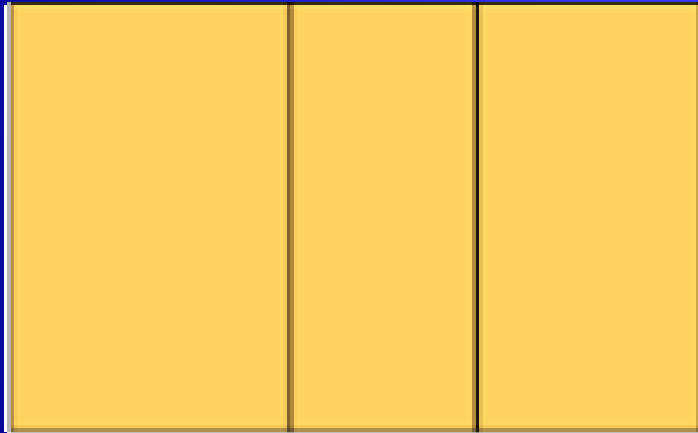# Predictive Methods

## *K. Gibert[1]*

*[1]Department of Statistics and Operation Research*

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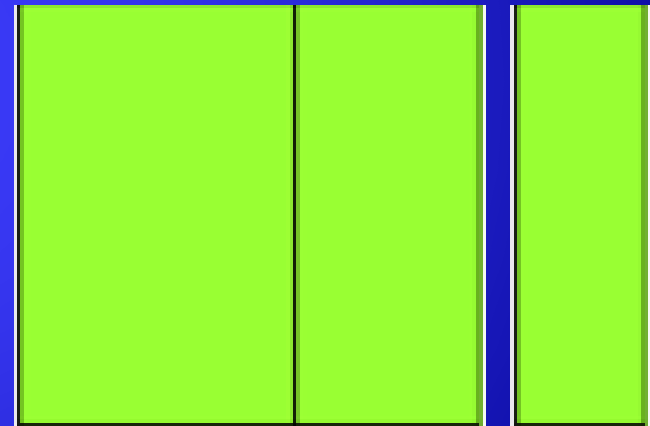
# Modelling



Cognition

Re-Cognition

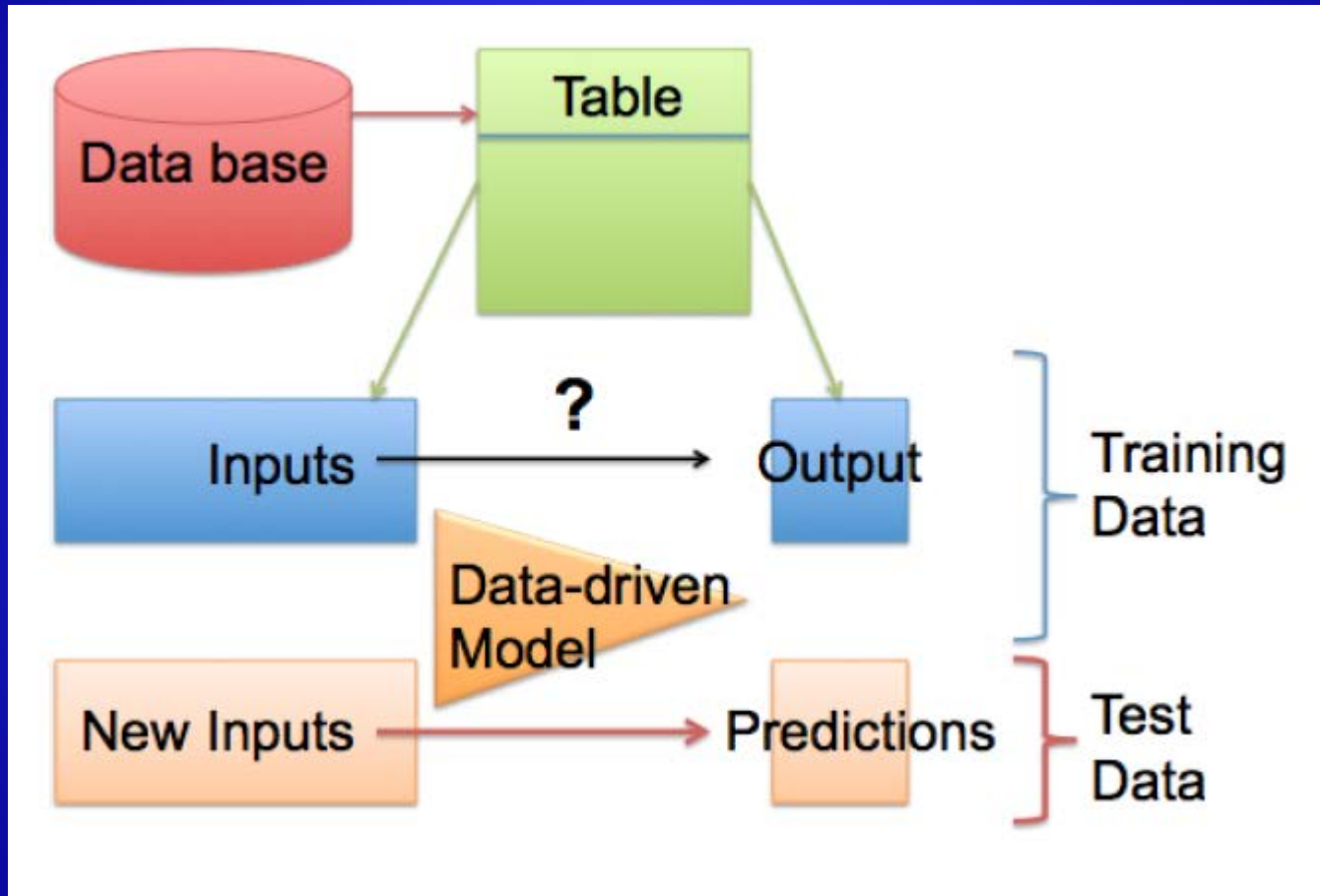**Socio-econ. Opinions Products**

**Inputs**   **Output(s)**

Data to explore

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

## Data= Fit+Error

- Fit:
  - Structural
  - Law governing the phenomenon
  - Analytic Function

- Error:
  - Random
  - Variability arround Fit (null expectation)
  - Probabilistic model

# Statistical models

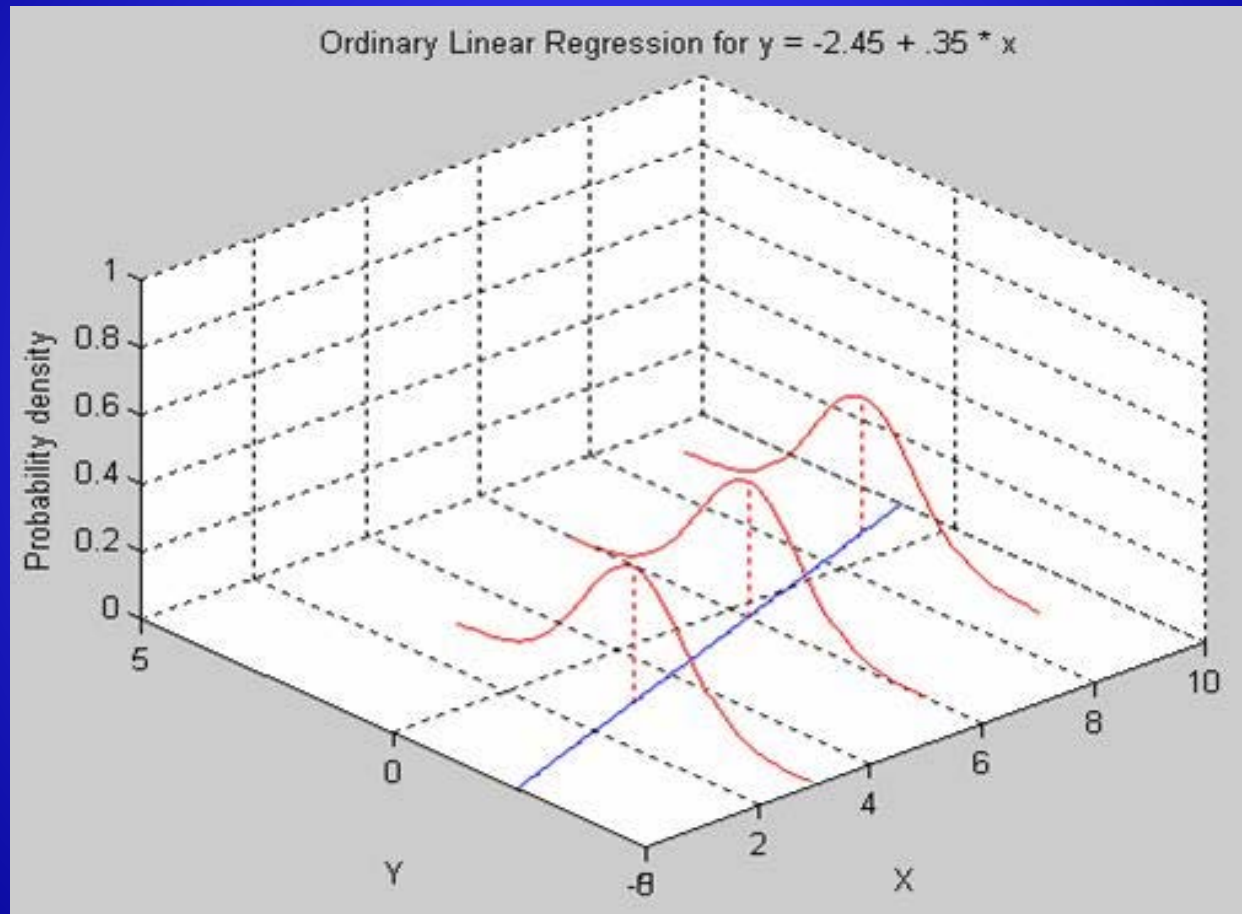- Determine the family of fits:
  - Linear
  - Quadratic
  - Exponential
  - .....

- Determine the law of error:
  - Normal
  - Poisson
  - Binomial....

# E1 Normally Distributed Error



Ordinary Linear Regression for y = -2.45 + .35 * x

# Linear Multiple regression

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

- *Formalization: I=i:n observations*
  *Y: Response variable*
  $X_{1 \dots} X_K$ : *ExplanatoryVariables*

  Find $\beta_0 \dots \beta_K$ such that

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

- *Assumptions:*

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

    ***Population regression line***

  - *Normality:* $\varepsilon_1, \dots, \varepsilon_n \sim \mathcal{N}(0, \sigma_i)$, $i = 1:n$

  - *Homokedasticity:* $Var[\varepsilon_i] = \sigma^2$ forall i

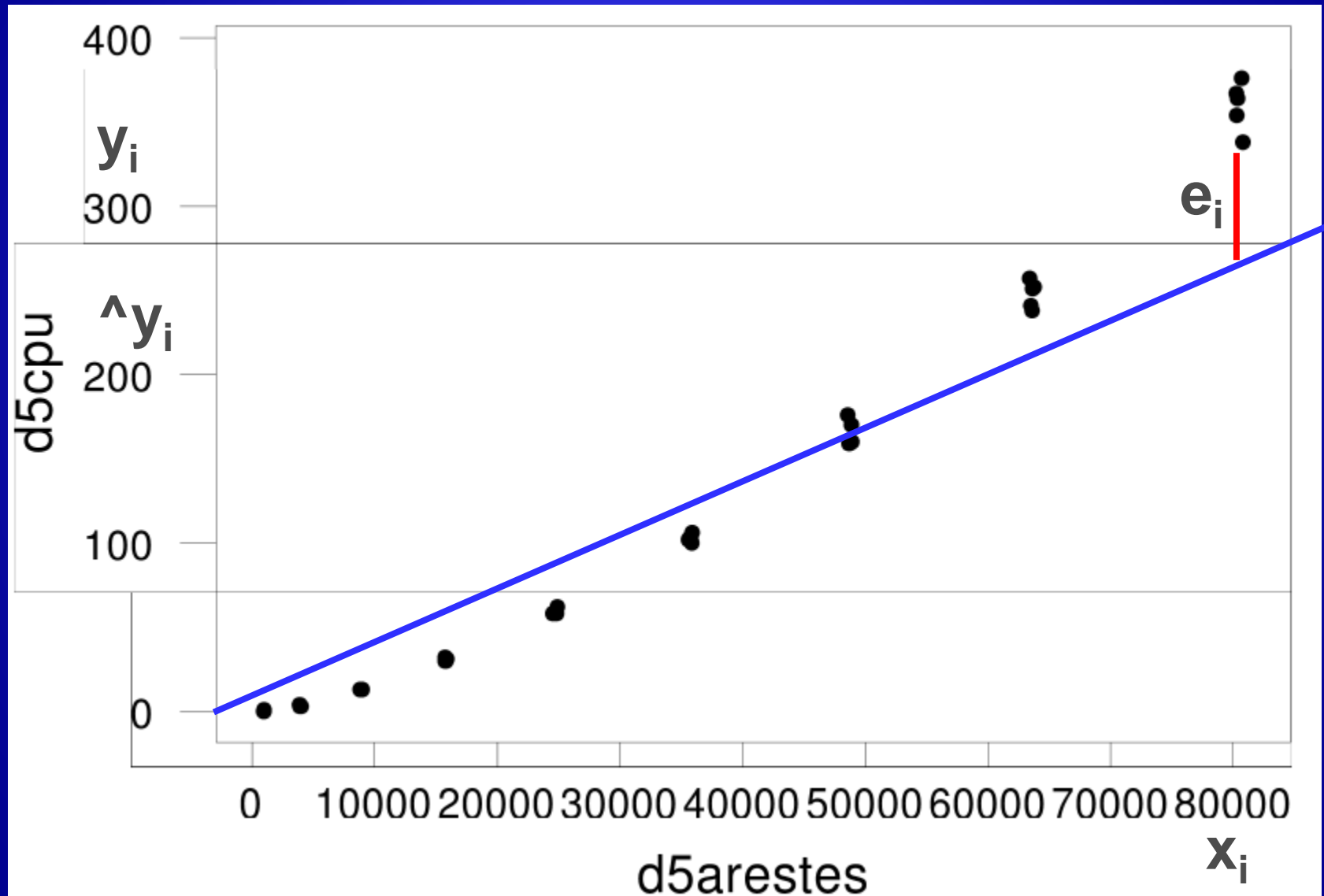  - *Independence:* $Cov(\varepsilon_i, \varepsilon_j) = 0$ forall i,,j

# What is "Linear"?

- Remember this:
- $Y = \beta_0 + \beta_1 X$

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

©R. Gibert

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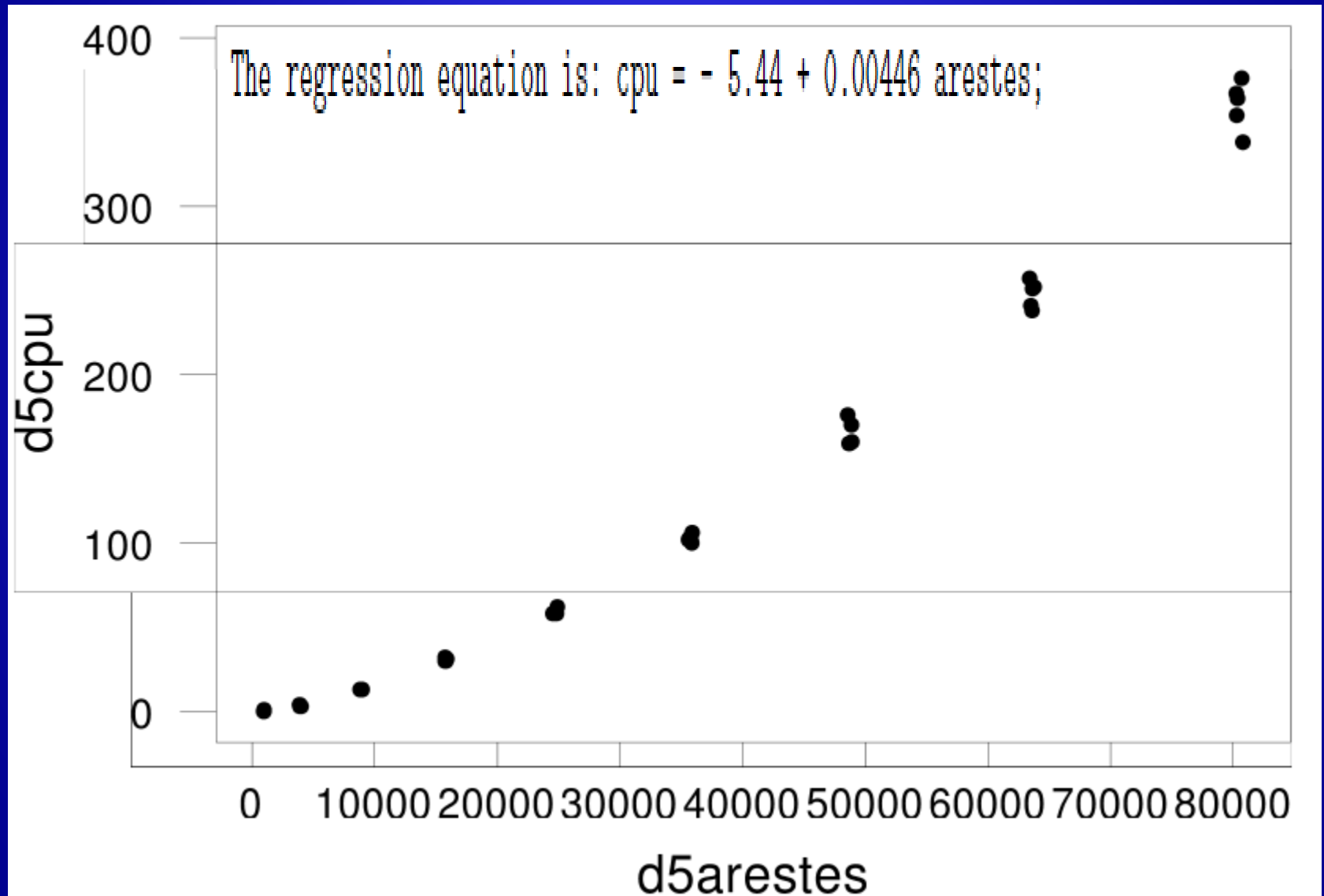**–** treatment algorithm vs graph size

# Minimum Least Squares solution

Find $\hat{\beta}_0, \hat{\beta}_1$ such that $\min_{\hat{\beta}_0, \hat{\beta}_1} \sum_{\forall i} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$
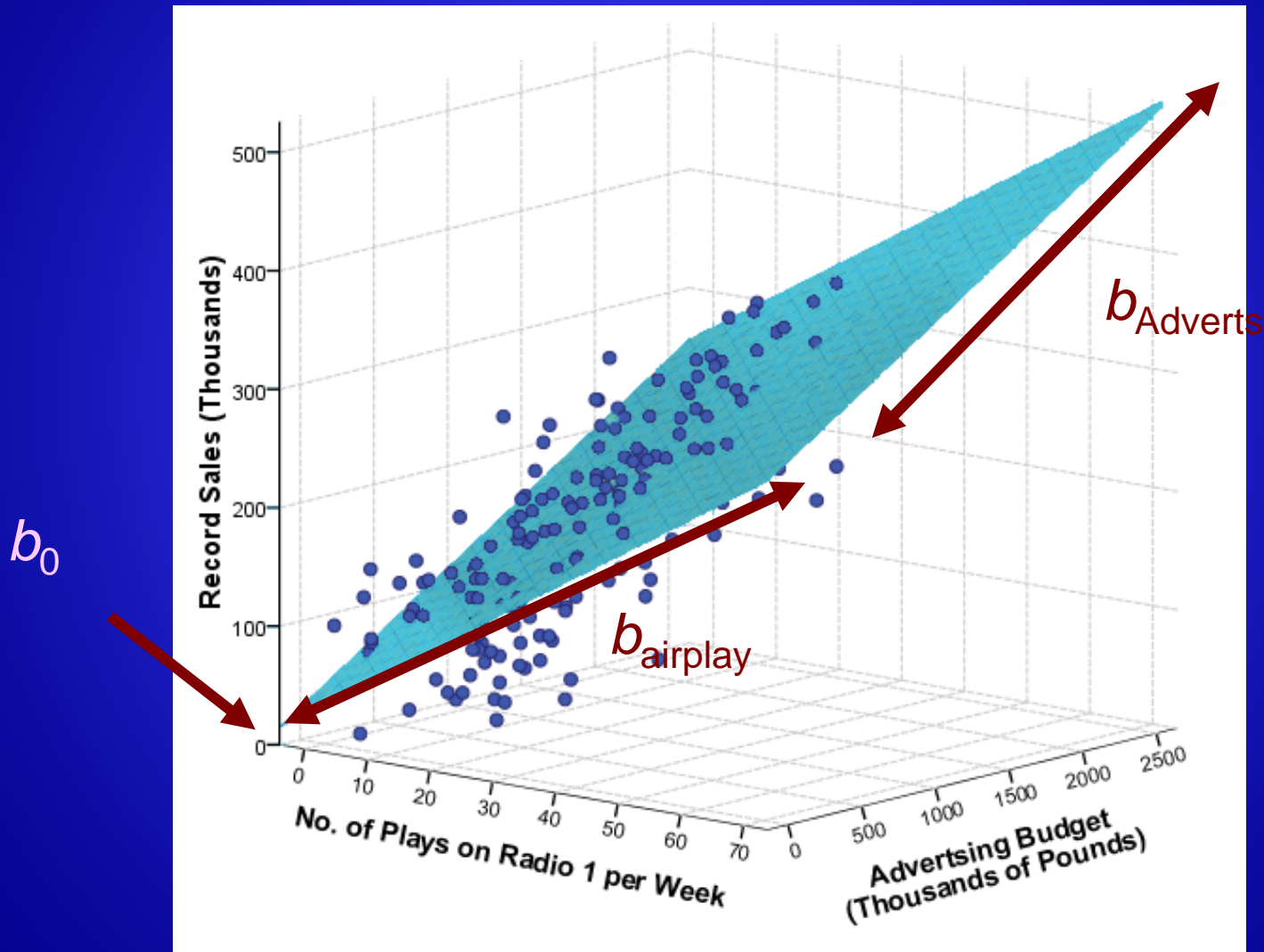
$$\hat{\beta}_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}$$

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

UPC

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



The regression equation is: cpu = − 5.44 + 0.00446 arestes;

# The Model with Two Predictors

# Matricial formulation

Regression fit criterion: $\min\limits_{r} E\left[\left(y_i - r(x_{i1}, \cdots, x_{ip})\right)^2\right]$

$$r(x_{i1}, \cdots, x_{ip}) = E\left[y_i \mid x_{i1}, \cdots, x_{ip}\right]$$

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

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

**Estimation of coefficients**

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

**In matrix notation**

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

# Geometric interpretation

$$y_i = \hat{y}_i + e$$

$$
\begin{matrix}
y & \hat{y} & e
\end{matrix}
$$

$$
\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} + \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}
$$

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

Criterion: $\quad \min\limits_{b_0,\dots,b_p} \sum\limits_{i=1}^{n} (e_i)^2 = \|e\|^2$

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

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

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

# Validation

- Technical Assumptions
  - normality, linearity, independence, homokedasticity

  - Tools
    - Graphical residuals analysis
    - Influence-point indicators (hi)

- Quality:
  - R2 (determination coeficient): 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 = SSexplainedByModel + SSerrror

Dividing all sides by SStotal:

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

# Quantify Goodnes of model

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

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

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

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 \, x_{1i} + \hat{\beta}_2 \, x_{2i} + .. + \hat{\beta}_k \, x_{Ki}$$

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

# The residuals

# Quantify Goodnes of model

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

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

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

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 \, x_{1i} + \hat{\beta}_2 \, x_{2i} + .. + \hat{\beta}_k \, x_{Ki}$$

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

UPC

# Quantify Goodnes of model

R$^2$ = proportion of explained variance

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

*assume linearity*

$$0 < R^2 < 1$$

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

For simple linear regression R2=Corr(Y, X)^2

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



The regression equation is: cpu = – 5.44 + 0.00446 arestes;

$S = 3.607$   $R-Sq = 97.8\%$

Floid algorithm

©R. Gibert

# Model inference

To test significance of the model

$$F = \frac{SSExplainedByModel}{SSError} \sim F_{(k-1,n-k)}$$

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

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

To test significance of a model term

**Both assume normality**

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



The regression equation is: cpu = - 5.44 + 0.00446 arestes;

S = 3.607   R-Sq = 97.8%

Numbers not enough

©R. Gibert

# Graphical residuals analysis



©K. Gibert

# Regression

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% |

4/28/2017

©*K. Gibert*

# Regression

# Regression

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

©K. Gibert

**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
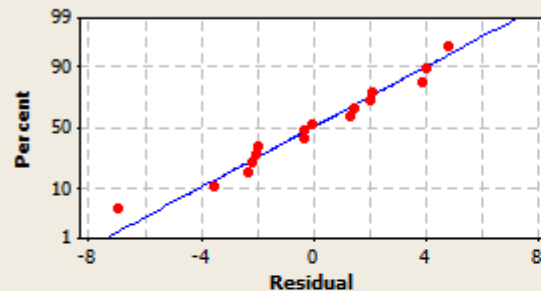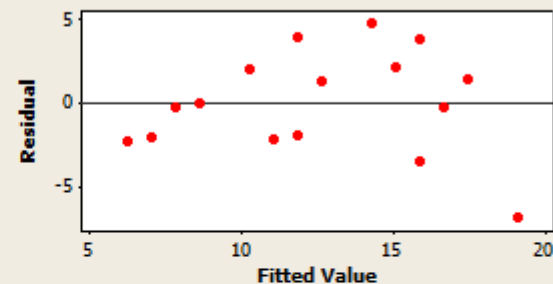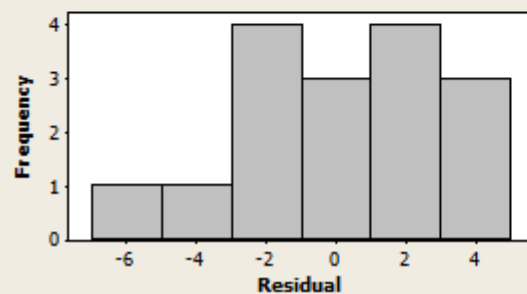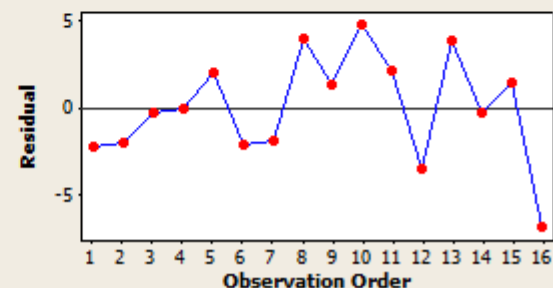
Versus Order

4/28/2017

**Fitted Line Plot**
y3 = 0,519 + 0,8087 x3

| S | 3,22542 |
| R-Sq | 61,7% |
| R-Sq(adj) | 59,0% |

**Residual Plots for y3**

Normal Probability Plot

Versus Fits

Histogram

Versus Order
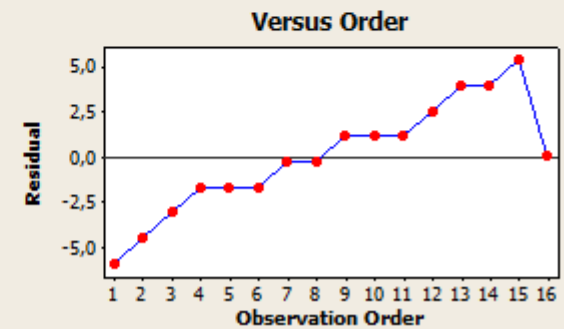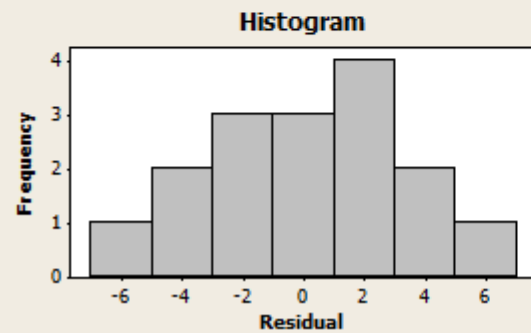
# Going further

- 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