

STATISTICAL DATA ANALYSIS IN EXCEL

Lecture 5

Linear Regression

dr. Petr Nazarov

petr.nazarov@crp-sante.lu

31-10-2011

Statistical data analysis in Excel.

5. Linear regression







Introduction

- covariation and correlation measures
- dependent and independent random variables
- scatter plot and linear trendline
- linear model

Testing for significance

- estimation of the noise variance
- interval estimations
- testing hypothesis about significance

Regression Analysis

- confidence and prediction
- multiple linear regression
- nonlinear regression



Dependent and Independent Variables





INTRODUCTION



Dependent and Independent Variables



INTRODUCTION



Covariance

Measure of Association between 2 Variables

indicate a positive relationship; negative values indicate a negative relationship. population sample $\sigma_{xy} = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{\sum (x_i - \mu_y)}$ $(x_i - m_x)$ $s_{xy} =$ n-1In Excel use function: 60 =COVAR(data) mice.xls 50 40 Ending weight 20 $s_{xy} = 39.8$ Ending weight VS. 10 Starting weight 0 10 40 hard to 20 30 50 0 Starting weight interpret

A measure of linear association between two variables. Positive values

Statistical data analysis in Excel.





Measure of Association between 2 Variables

Correlation (Pearson product moment correlation coefficient)

A measure of linear association between two variables that takes on values between -1 and +1. Values near +1 indicate a strong positive linear relationship, values near -1 indicate a strong negative linear relationship; and values near zero indicate the lack of a linear relationship.

population

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} = \frac{\sum (x_i - m_x)(y_i - m_y)}{\sigma_x \sigma_y N}$$



sample

$$r_{xy} = \frac{s_{xy}}{s_x s_y} = \frac{\sum (x_i - m_x)(y_i - m_y)}{s_x s_y (n-1)}$$

In Excel use function:

$$r_{xy} = 0.94$$

Statistical data analysis in Excel.





Experiments

Cells are grown under different temperature conditions from 20° to 40°. A researched would like to find a dependency between T and cell number.

Dependent variable

The variable that is being predicted or explained. It is denoted by y.

Independent variable

The variable that is doing the predicting or explaining. It is denoted by **x**.

45

X



Experiments

Simple linear regression

Regression analysis involving one independent variable and one dependent variable in which the relationship between the variables is approximated by a straight line.

Building a regression means finding and tuning the model to explain the behaviour of the data





Experiments

Regression model

The equation describing how y is related to x and an error term; in simple linear regression, the regression model is $y = \beta_0 + \beta_1 x + \varepsilon$

Regression equation

The equation that describes how the mean or expected value of the dependent variable is related to the independent variable; in simple linear regression, $E(y) = \beta_0 + \beta_1 x$



Model for a simple linear regression:

$$y(x) = \beta_1 x + \beta_0 + \varepsilon$$



Regression Model and Regression Line

$$y(x) = \beta_1 x + \beta_0 + \varepsilon$$

Cr

CENTRE DE RECHERCHE PUBLIC





Estimated regression equation

The estimate of the regression equation developed from sample data by using the least squares method. For simple linear regression, the estimated regression equation is $y = b_0 + b_1 x$

$$\hat{y}(x) = b_1 x + b_0$$

 $E[y(x)] = b_1 x + b_0$

 $y(x) = \beta_1 x + \beta_0 + \varepsilon$

cells.xls

CENTRE DE RECHERCHE PUBLI

1. Make a scatter plot for the data.



2. Right click to "Add Trendline". Show equation.





SIMPLE LINEAR REGRESSION

Overview





Statistical data analysis in Excel. 5. Linear regression





Experiments

Least squares method

A procedure used to develop the estimated regression equation.

The objective is to minimize

 $\sum (y_i - \hat{y}_i)^2$

 y_i = observed value of the dependent variable for the *i*th observation \hat{y}_i = estimated value of the dependent variable for the *i*th observation

Slope:

$$b_1 = \frac{\sum (x_i - m_x)(y_i - m_y)}{(x_1 - m_x)^2}$$
Intersect:

$$b_0 = m_y - b_1 m_x$$

Coefficient of Determination



The Main Equation

$$SST = SSR + SSE$$





ANOVA and Regression: Testing for Significance



Coefficient of Determination





Statistical data analysis in Excel. 5. Linear regression

Correlation coefficient

A measure of the strength of the linear relationship between two variables (previously discussed in Lecture 1).





$$SSE = \sum (y_i - \mathbf{y}_i)^2$$
$$SST = \sum (y_i - \overline{y})^2$$

$$SSR = \sum \left(\oint_{i} - \overline{y} \right)^{2}$$

Coefficient of determination

0

A measure of the goodness of fit of the estimated regression equation. It can be interpreted as the proportion of the variability in the dependent variable y that is explained by the estimated regression equation.

SST = SSR + SSE





Assumptions

Assumptions for Simple Linear Regression

- **1.** The error term ε is a random variable with 0 mean, i.e. $E[\varepsilon]=0$
- **2.** The variance of $\boldsymbol{\varepsilon}$, denoted by $\boldsymbol{\sigma}^2$, is the same for all values of \boldsymbol{x}
- **3.** The values of $\boldsymbol{\varepsilon}$ are independent
- **3.** The term $\boldsymbol{\varepsilon}$ is a normally distributed variable





Statistical data analysis



Confidence and Prediction

Confidence interval

The interval estimate of the mean value of y for a given value of x.

Prediction interval

The interval estimate of an individual value of y for a given value of x.





Residuals



C

CENTRE DE RECHERCHE PUBLIC



TESTING FOR SIGNIFICANCE

Sampling Distribution for *b*₁

If assumptions for ϵ are fulfilled, then the sampling distribution for b_1 is as follows:

$$y(x) = \beta_1 x + \beta_0 + \varepsilon$$
$$\hat{y}(x) = b_1 x + b_0$$



Interval Estimation for β_1

$$\beta_1 = b_1 \pm t_{\alpha/2}^{(n-2)} \frac{\sigma}{\sqrt{\sum (x_i - m_x)^2}}$$

$$\beta_1 = b_1 \pm t_{\alpha/2}^{(n-2)} SE$$

Statistical data analysis in Excel.



TESTING FOR SIGNIFICANCE

Test for Significance

$$H_0$$
: β₁ = 0 insignificant
 H_a : β₁ ≠ 0

1. Build a t-test statistics.

$$t = \frac{b_1}{\sigma_{b_1}} = \frac{b_1}{s} \sqrt{\sum (x_i - m_x)^2}$$



2. Calculate p-value for *t*

p-value approach:Reject H_0 if p-value $\leq \alpha$ Critical value approach:Reject H_0 if $t \leq -t_{a/2}$ or if $t \geq t_{a/2}$

where $t_{a/2}$ is based on a t distribution with n - 2 degrees of freedom.





Example

cells.xls

- **1.** Calculate manually b_1 and b_0
- Intercept b0= -191.008119 Slope b1= 15.3385723
- In Excel use the function:
- = INTERCEPT(y, x)
- \Rightarrow = SLOPE(y,x)

2. Let's do it automatically Tools \rightarrow Data Analysis \rightarrow Regression

SUMMARY OUTPUT

Regression Statistics						
Multiple R	0.950842308					
R Square	0.904101095					
Adjusted R Square	0.899053784					
Standard Error	31.80180903					
Observations	21					

ANOVA

	df		SS	MS	F	Significance F	
Regression		1	181159.2853	181159.3	179.1253	4.01609E-11	
Residual		19	19215.7461	1011.355			
Total		20	200375.0314				

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	-191.0081194	35.07510626	-5.445689	2.97E-05	-264.4211603	-117.5950784	-264.4211603	-117.5950784
X Variable 1	15.33857226	1.146057646	13.38377	4.02E-11	12.93984605	17.73729848	12.93984605	17.73729848

Statistical data analysis in Excel.



rana.txt

A biology student wishes to determine the relationship between temperature and heart rate in leopard frog, *Rana pipiens*. He manipulates the temperature in 2° increment ranging from 2 to 18°C and records the heart rate at each interval. His data are presented in table rana.txt

- 1) Build the model and provide the p-value for linear dependency
- 2) Provide interval estimation for the slope of the dependency
- 3) Estmate 95% prediction interval for heart rate at 15°

Multiple Regression







Multiple Regression





Non-Linear Regression

FIGURE 15.12 LOGISTIC REGRESSION EQUATION FOR $\beta_0 = -7$ AND $\beta_1 = 3$



$$E(y) = P(y = 1 | x_1, x_2, ..., x_p) = \frac{\exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_p x_p)}{1 + \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_p x_p)}$$



CASE STUDY

Correlation Analysis of Transcriptomic Data

Gene regulatory networks (GRN) in living cells can be considered as extremely complex information processing systems. Despite their complexity, the main feature of the GRN is their robustness and ability to form a proper biochemical respond to a wide range of extracellular conditions. The knowledge about the part of GRN related to a specific bio-function of cellular process is of extreme importance for controlling them. Another important aspect of understanding cell functionality is linked to knowledge about the regulatory effect of small non-coding micro-RNA (miRNA). miRNAs influence most fundamental biological processes by ultimately altering the expression levels of proteins either through degradation of mRNA or through interference with mRNA translation. miRNAs tend to have long half lives and therefore represent promising candidates to be used as disease markers and therapeutic targets.

Being a reverse-engineering task, the GRN reconstruction is highly challenging, and requires analysis of large sets of experimental data. One of the straightest ways to reconstruct GRN is based on co-expression (CE) analysis of transcriptomic data from cDNA microarrays. Two significantly co-expressed genes or a gene and miRNA have the same or inverted expression profile over a number of samples. Biologically this is a good evidence for either a direct interaction between the genes or their mutual participation in the same biological function.

The performance of the software was tested using public mRNA and miRNA expression data from 14 various cell lines (A498, ACHN, CAKI1, CCRFCEM, HCT15, HL60, K562, MALME3M, MCF7, MOLT4, NCIH226, NCIH522, RPMI8226, SKOV3). Data from 42 Affymetrix® HGU133plus2 arrays and 14 miRNA custom microarray experiments were downloaded from public repositories (ref. E-MTAB-37 and E-MEXP-1029, <u>http://www.ebi.ac.uk</u>), normalized and analyzed.

Tool: http://edu.sablab.net/biostat2/coexpress.zip

Data: http://edu.sablab.net/biostat2/data-mir-mrna_14cl.zip





Thank you for your attention

to be continued...

