

STATISTICAL DATA ANALYSIS IN EXCEL

Lecture 5

Linear Regression

dr. Petr Nazarov

petr.nazarov@crp-sante.lu

14-01-2013

◆ **Introduction**

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

◆ **Testing for significance**

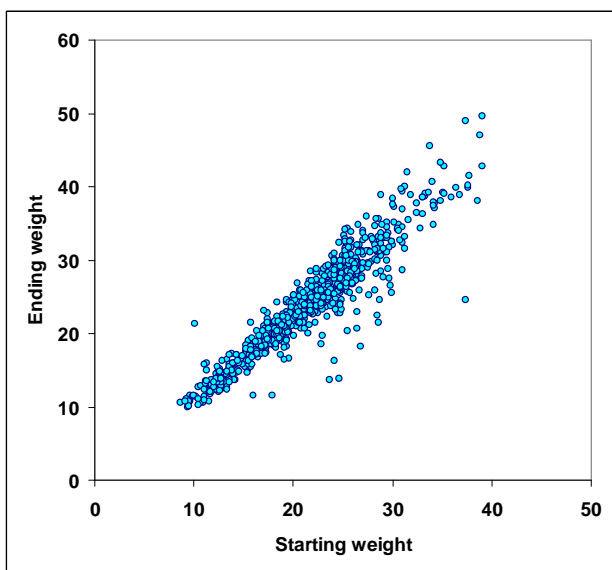
- ◆ estimation of the noise variance
- ◆ interval estimations
- ◆ testing hypothesis about significance

◆ **Regression Analysis**

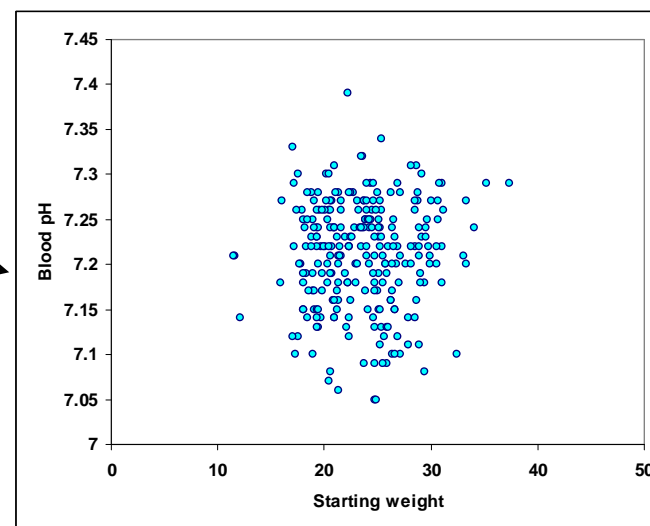
- ◆ confidence and prediction
- ◆ multiple linear regression
- ◆ nonlinear regression

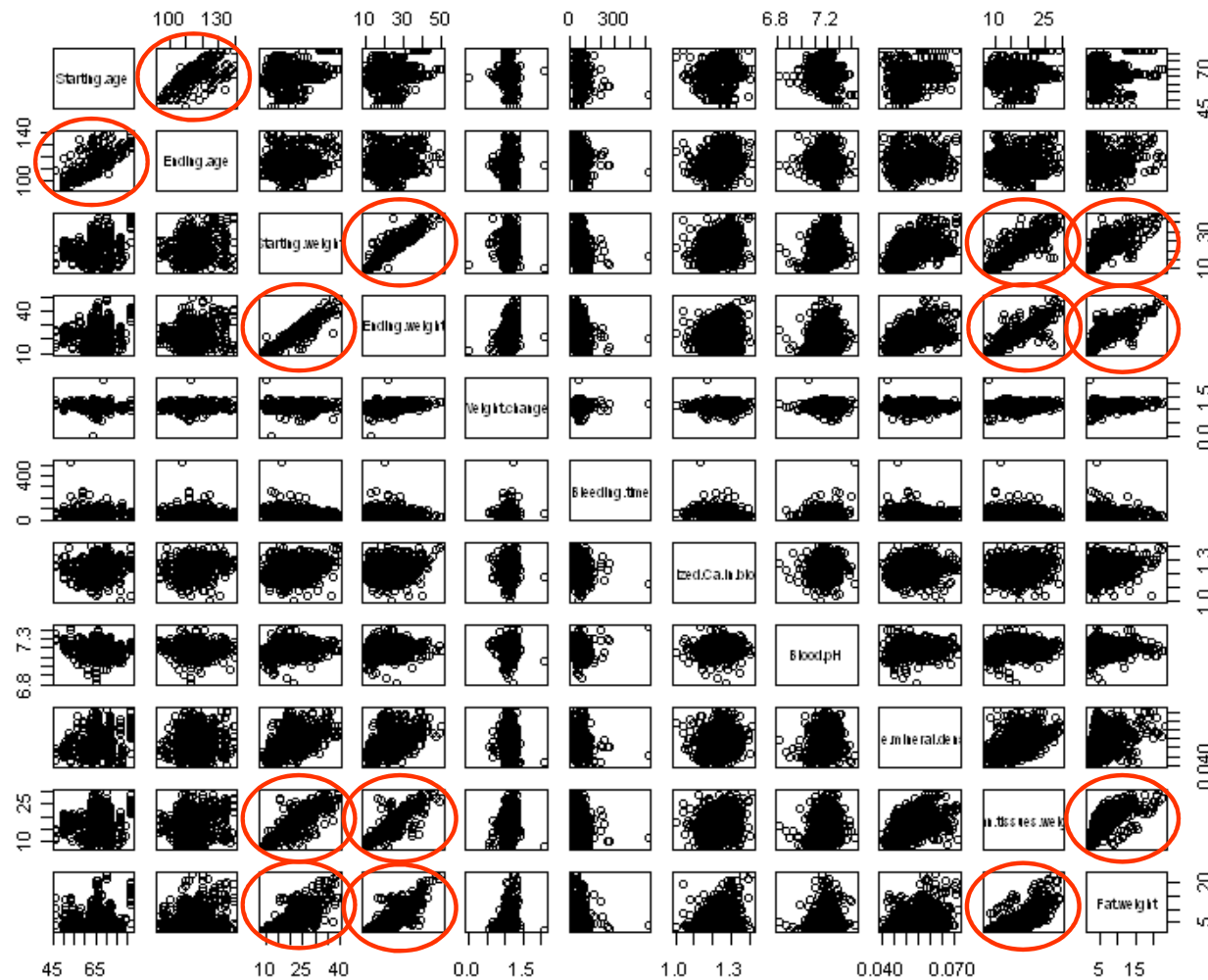
mice.xls

Ending weight vs. Starting weight



Blood pH vs. Starting weight





Measure of Association between 2 Variables

Covariance

A measure of linear association between two variables. Positive values indicate a positive relationship; negative values indicate a negative relationship.

population

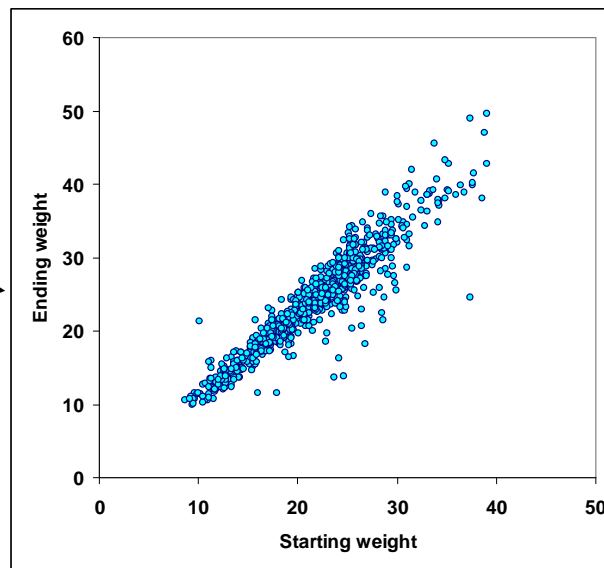
$$\sigma_{xy} = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{N}$$

sample

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

mice.xls

Ending weight vs.
Starting weight



In Excel use function:

◆ =COVAR(data)

$s_{xy} = 39.8$

hard to interpret

Measure of Association between 2 Variables

Correlation (Pearson product moment correlation coefficient)

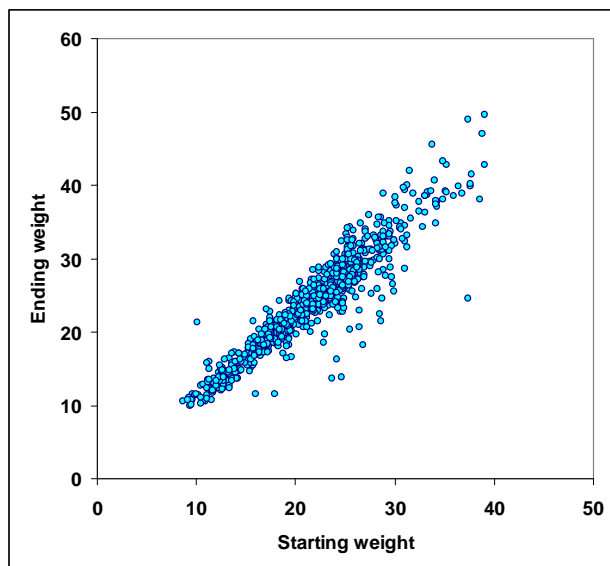
A measure of linear association between two variables that takes on values between -1 and +1. Values near +1 indicate a strong positive linear relationship, values near -1 indicate a strong negative linear relationship; and values near zero indicate the lack of a linear relationship.

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

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

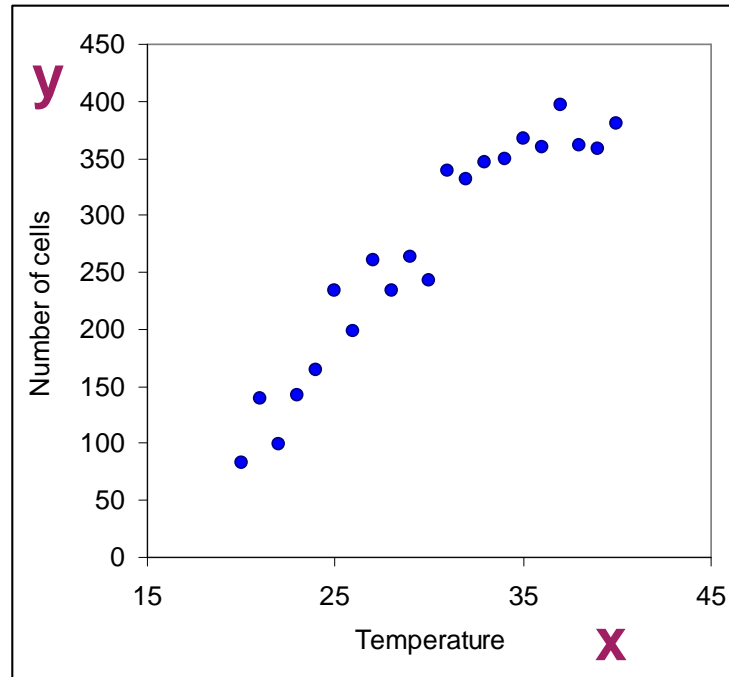


In Excel use function:

◆ =CORREL(data)

$$r_{xy} = 0.94$$

Temperature	Cell Number
20	83
21	139
22	99
23	143
24	164
25	233
26	198
27	261
28	235
29	264
30	243
31	339
32	331
33	346
34	350
35	368
36	360
37	397
38	361
39	358
40	381



Cells are grown under different temperature conditions from 20° to 40°. A researcher 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 .

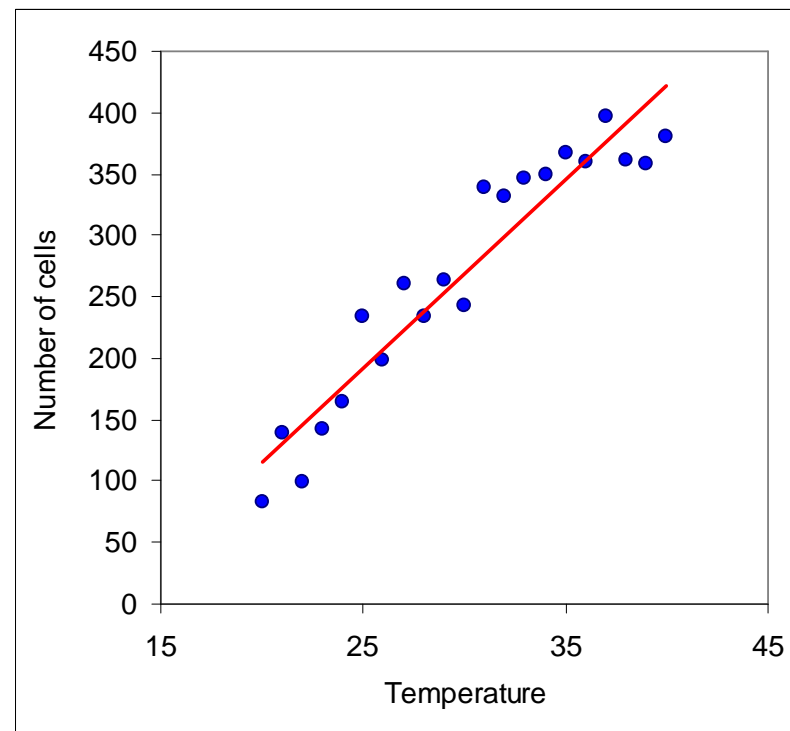
Independent variable

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

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*



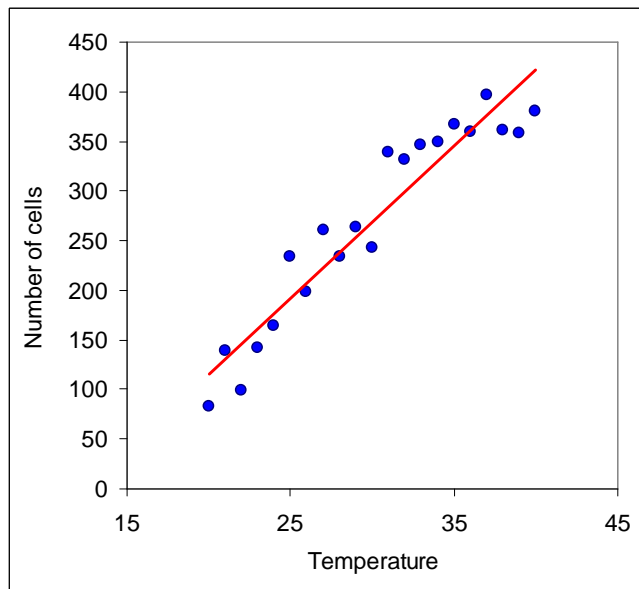
Regression model

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,

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

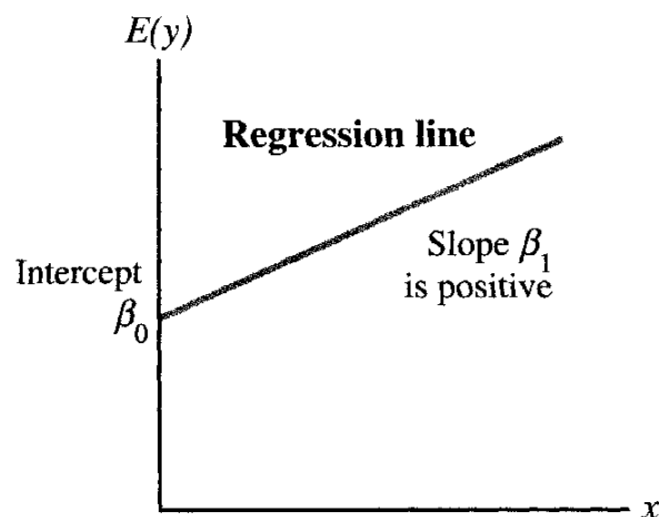


◆ Model for a simple linear regression:

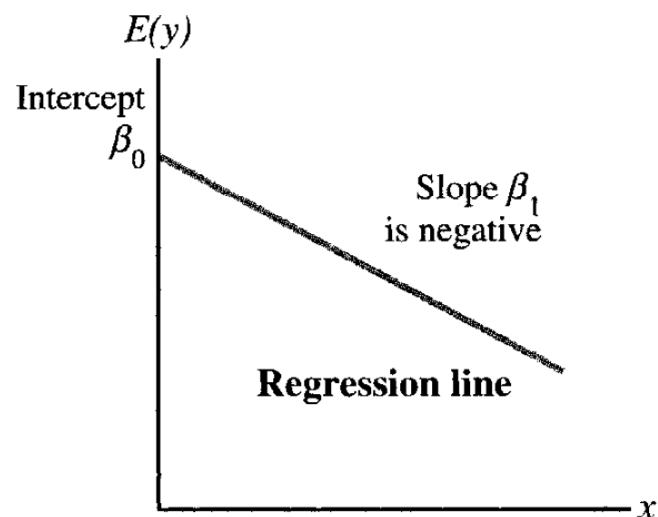
$$y(x) = \beta_1 x + \beta_0 + \varepsilon$$

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

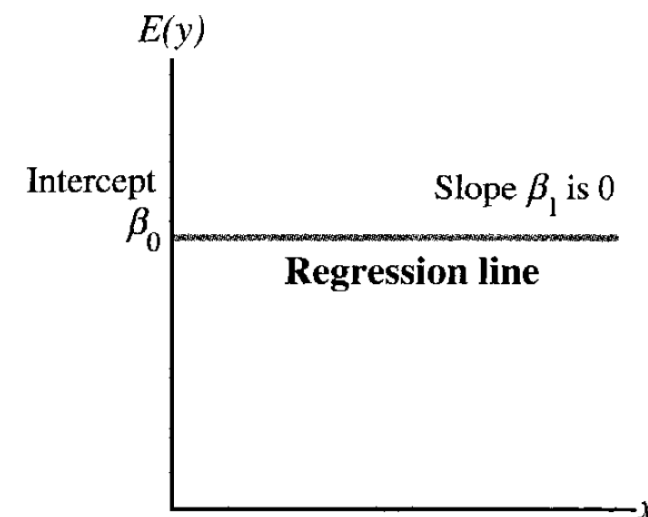
**Panel A:
Positive Linear Relationship**



**Panel B:
Negative Linear Relationship**



**Panel C:
No Relationship**



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 $y = b_0 + b_1x$

$$y(x) = \beta_1x + \beta_0 + \varepsilon$$

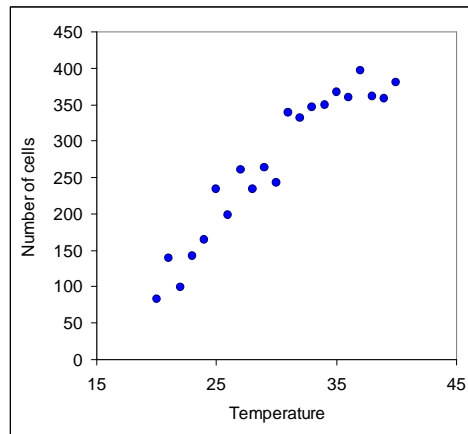


$$\hat{y}(x) = b_1x + b_0$$

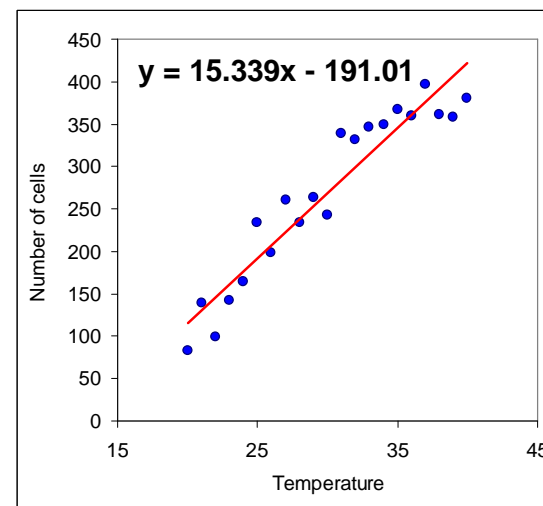
cells.xls

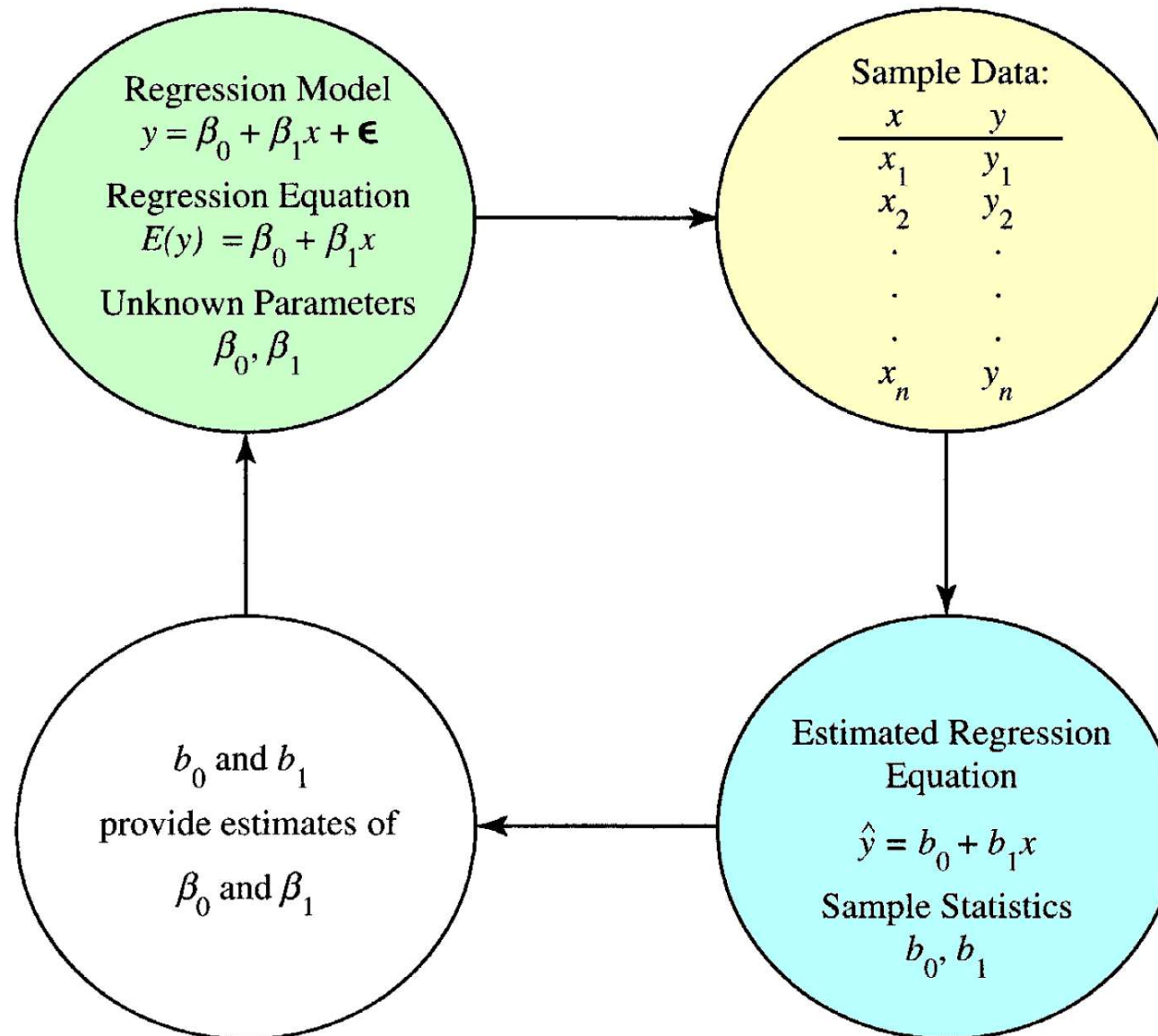
$$E[y(x)] = b_1x + b_0$$

1. Make a scatter plot for the data.



2. Right click to “Add Trendline”. Show equation.





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:

$$b_0 = m_y - b_1 m_x$$

Sum squares due to **error**

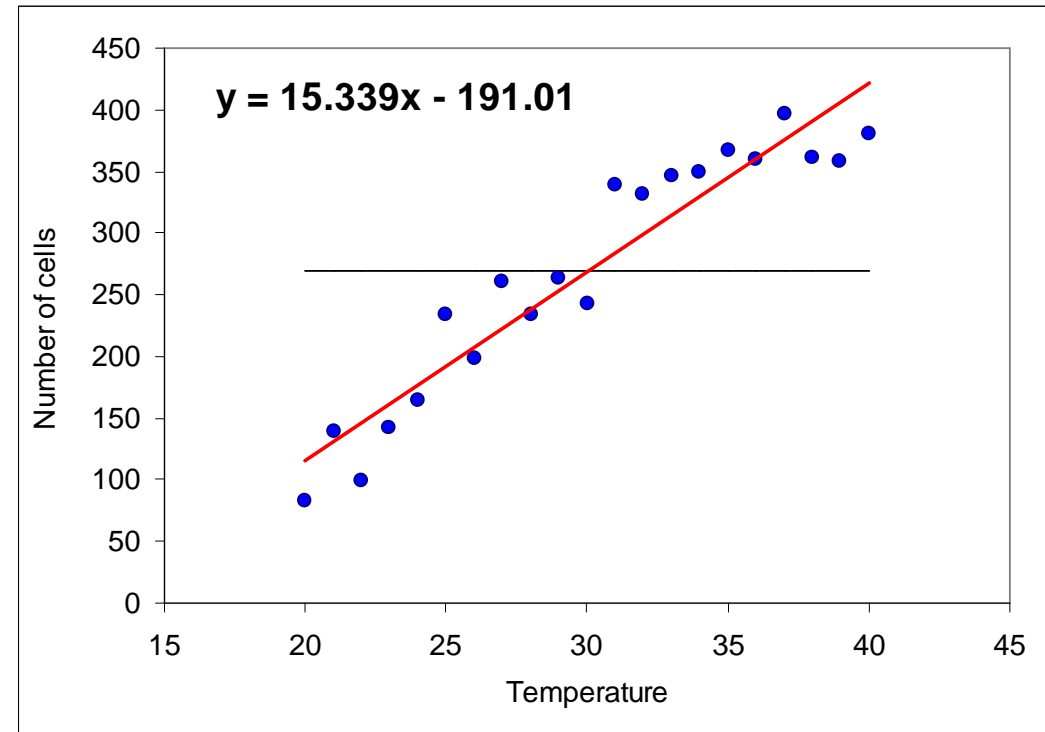
$$SSE = \sum (y_i - \hat{y}_i)^2$$

Sum squares **total**

$$SST = \sum (y_i - \bar{y})^2$$

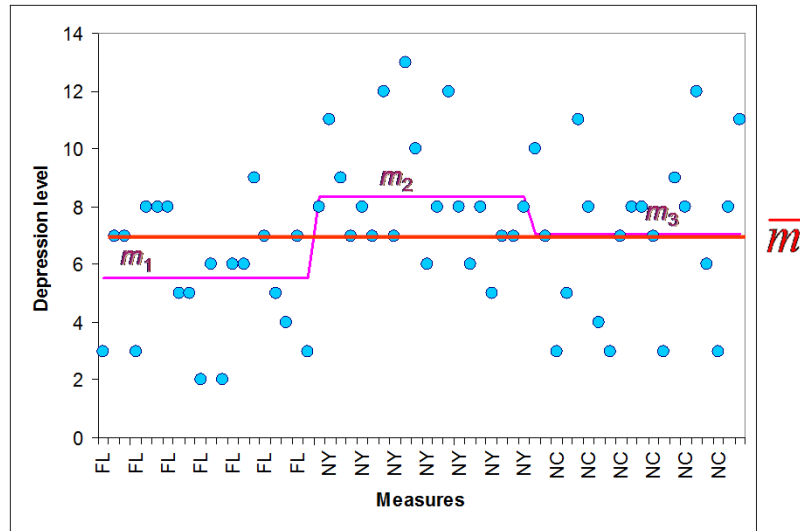
Sum squares **due to regression**

$$SSR = \sum (\hat{y}_i - \bar{y})^2$$

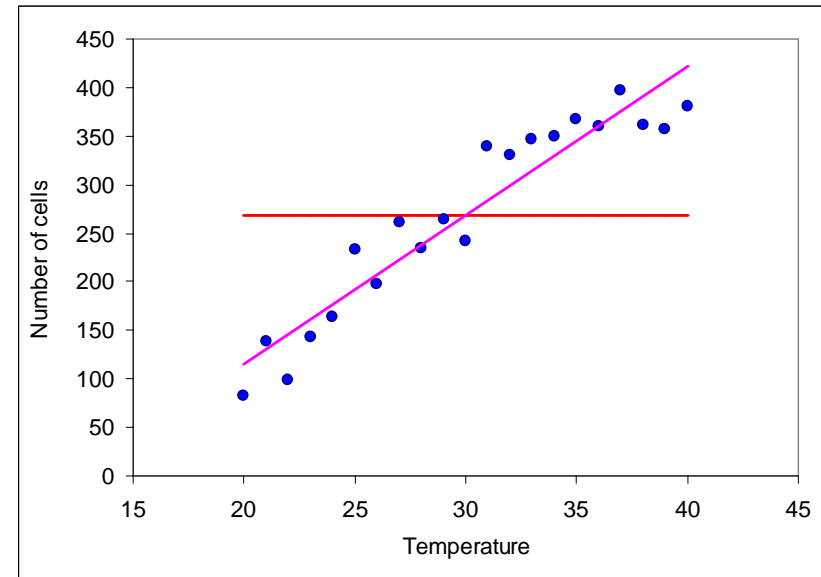


The Main Equation

$$SST = SSR + SSE$$



$$SST = SSTR + SSE$$



$$SST = SSR + SSE$$

$H_0: \beta_1 = 0$ *insignificant*

$H_a: \beta_1 \neq 0$

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

$$SST = \sum (y_i - \bar{y})^2$$

$$SSR = \sum (\hat{y}_i - \bar{y})^2$$

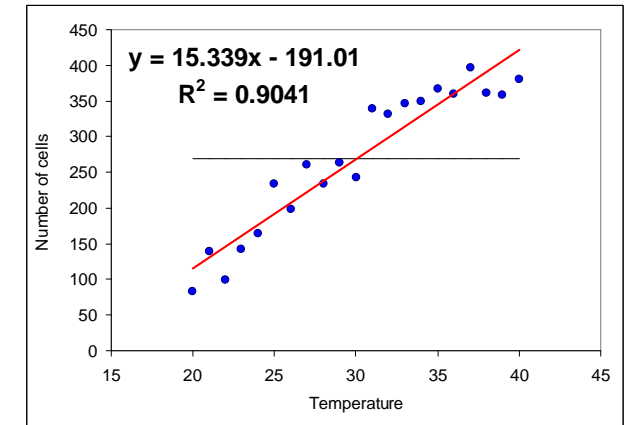
$$SST = SSR + SSE$$

Coefficient of determination

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.

Correlation coefficient

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



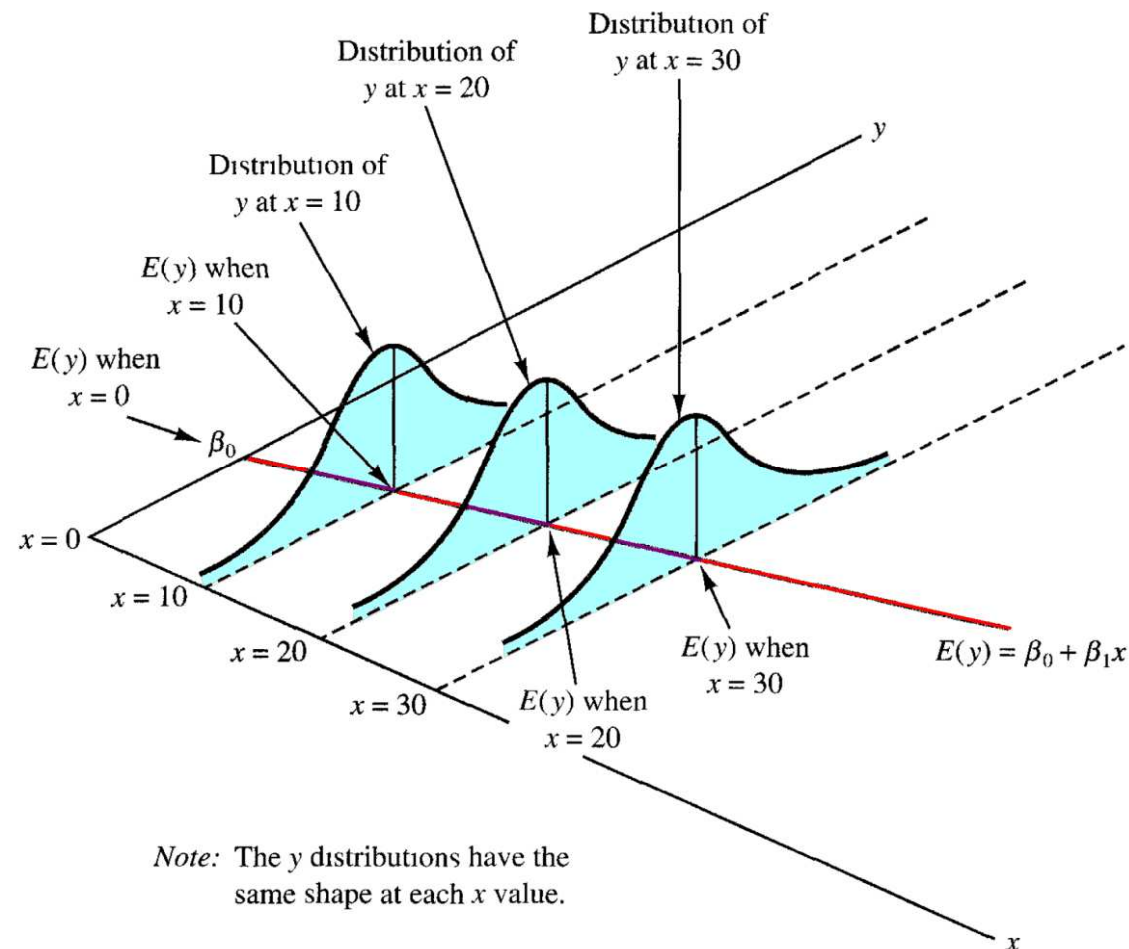
$$R^2 = \frac{SSR}{SST}$$

$$r = \text{sign}(b_1) \sqrt{R^2}$$

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 ε , denoted by σ^2 , is the same for all values of x
3. The values of ε are independent
3. The term ε is a normally distributed variable

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

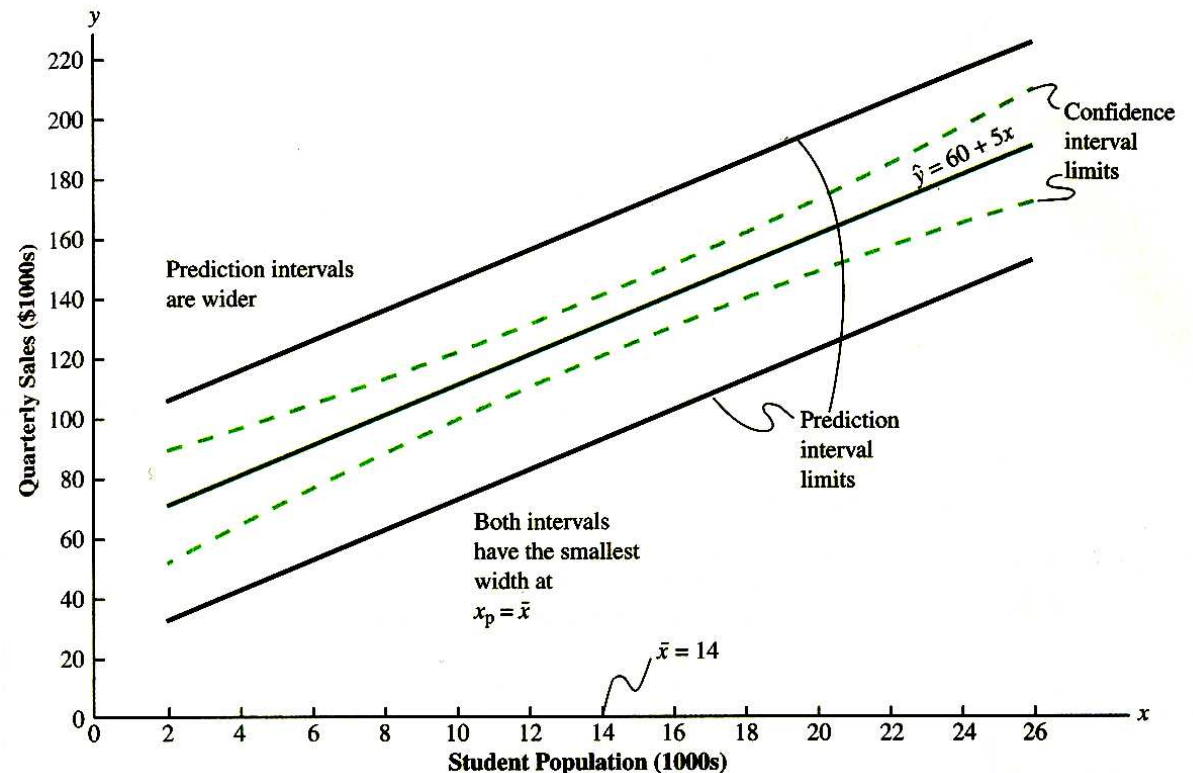
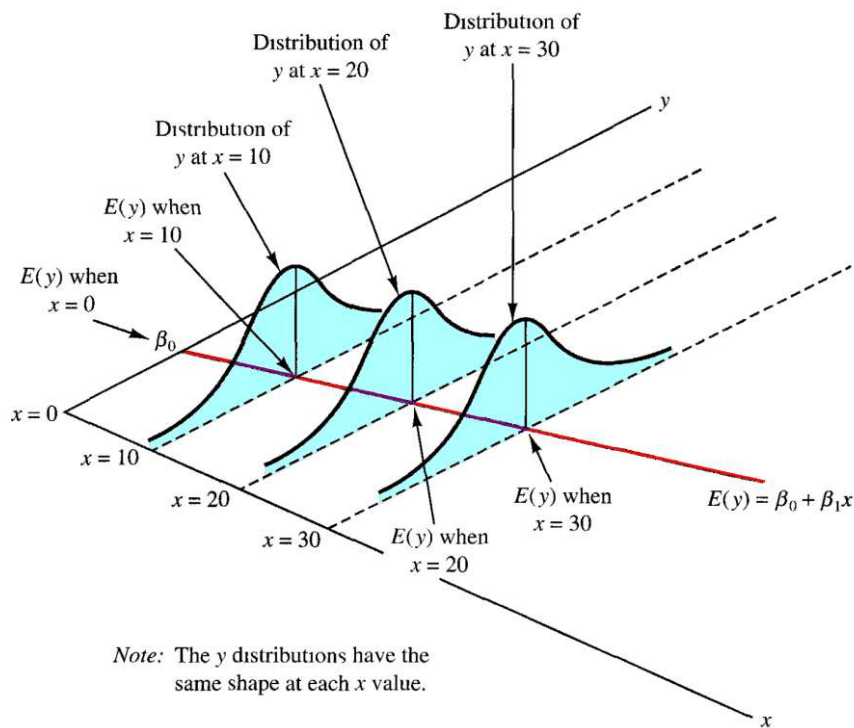


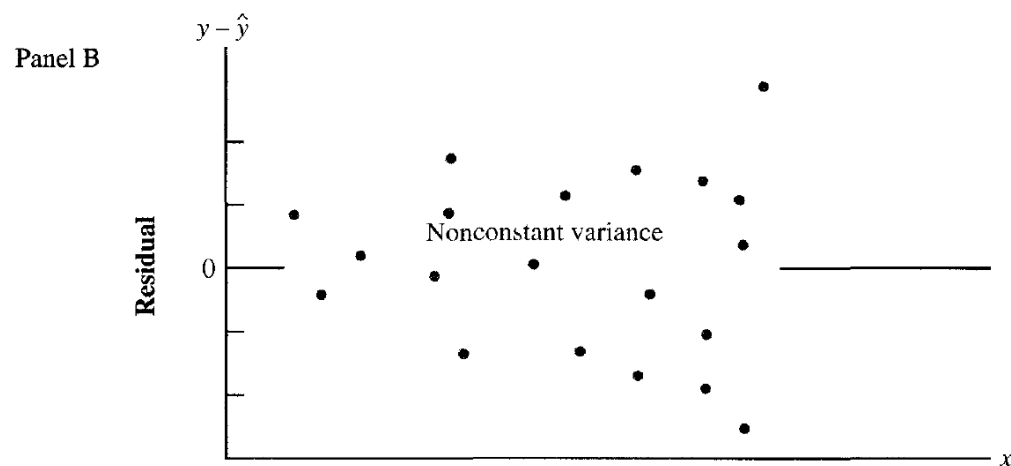
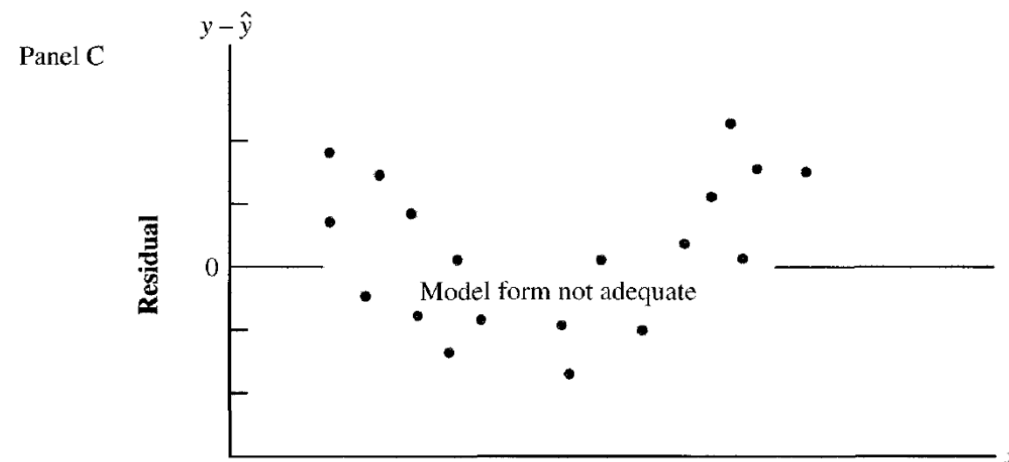
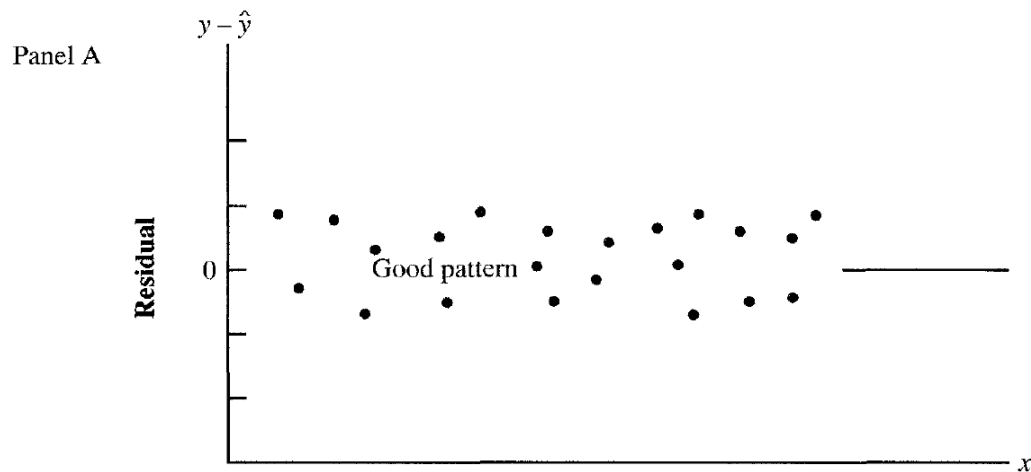
Confidence interval

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 .





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

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

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

Expected value

$$E[b_1] = \beta_1$$

Variance

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

= Standard Error

Distribution:

normal

Interval Estimation for β_1

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

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

$H_0: \beta_1 = 0$ *insignificant*

$H_a: \beta_1 \neq 0$

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 $p\text{-value} \leq \alpha$

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

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

1. Build a F-test statistics.

$$F = \frac{MSR}{MSE}$$

$$MSR = \frac{SSR}{\text{Number of independent variables}}$$

2. Calculate a p-value

cells.xls

1. Calculate manually b_1 and b_0

Intercept $b_0 = -191.008119$
 Slope $b_1 = 15.3385723$

In Excel use the function:

◆ = INTERCEPT(y, x)

◆ = SLOPE(y, x)

2. Let's do it automatically [Tools](#) → [Data Analysis](#) → [Regression](#)

SUMMARY OUTPUT

<i>Regression Statistics</i>	
Multiple R	0.950842308
R Square	0.904101095
Adjusted R Square	0.899053784
Standard Error	31.80180903
Observations	21

ANOVA

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

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>	<i>Lower 95.0%</i>	<i>Upper 95.0%</i>
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

rana.xls

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) Estimate 95% prediction interval for heart rate at 15°

Thank you for your attention

to be continued...

