



CENTRE DE RECHERCHE PUBLIC

PhD Course Advanced Biostatistics

Lecture 4 Linear Models: ANOVA and Linear Regression

dr. P. Nazarov

petr.nazarov@crp-sante.lu

16-12-2014

L4. Linear models







ANOVA (L4.1)

- 1-factor ANOVA
- Multifactor ANOVA
- ♦ Experimental design

Linear regression (L4.2)





Why ANOVA ?

Means for more than 2 populations We have measurements for 5 conditions. Are the means for these conditions equal?

Validation of the effects

We assume that we have several factors affecting our data. Which factors are most significant? Which can be neglected? If we would use pairwise comparisons, what will be the probability of getting error?

Number of comparisons:

$$=\frac{5!}{2!3!}=10$$

 C_{2}^{5}

Probability of an error: $1-(0.95)^{10} = 0.4$





http://easylink.playstream.com/affymetrix/ambsymposium/partek_08.wvx

L4. Linear models





Example

As part of a long-term study of individuals 65 years of age or older, sociologists and physicians at the Wentworth Medical Center in upstate New York investigated the relationship between geographic location and depression. A sample of 60 individuals, all in reasonably good health, was selected; 20 individuals were residents of Florida, 20 were residents of New York, and 20 were residents of North Carolina. Each of the individuals sampled was given a standardized test to measure depression. The data collected follow; higher test scores indicate higher levels of depression.

Q: Is the depression level same in all 3 locations?

depression.txt

1. Good health respondents								
Florida	New York	N. Carolina						
3	8	10						
7	11	7						
7	9	3						
3	7	5						
8	8	11						
8	7	8						

$$H_0: \mu_1 = \mu_2 = \mu_3$$

 $H_a:$ not all 3 means are equal







Meaning

$$H_0: \mu_1 = \mu_2 = \mu_3$$

 $H_{\rm a}$: not all 3 means are equal









Assumption for ANOVA

Assumptions for Analysis of Variance

1. For each population, the response variable is normally distributed

2. The variance of the respond variable, denoted as σ^2 is the same for all of the populations.

3. The observations must be independent.









Some Calculations

Parameter	Florida	New York	N. Carolina
m=	5.55	8.35	7.05
overall mean=	6.98333		
var=	4.5763	4.7658	8.0500

Let's estimate the variance of sampling distribution. If H_0 is true, then all m_i belong to the same distribution

$$\sigma_m^2 = \frac{\sigma^2}{n}$$

$$m_3 \quad \mu \quad m_1 \quad m_2$$

$$\sigma_m^2 = \frac{\sum_{i=1}^{k} (m_i - \overline{m})^2}{k - 1} = \frac{(5.55 - 6.98)^2 + (8.35 - 6.98)^2 + (7.05 - 6.98)^2}{3 - 1} = 1.96$$

$$\sigma^2 = n\sigma_m^2 = 20 \times 1.96 = 39.27$$
 - this is called between-treatment estimate, works only at H₀

At the same time, we can estimate the variance just by averaging out variances for each populations:

$$\sigma^2 = \frac{\sum_{i=1}^k \sigma_i}{k} = \frac{4.58 + 4.77 + 8.05}{3} = 5.8$$

- this is called within-treatment estimate

Does between-treatment estimate and within-treatment estimate give variances of the same "population"?

L4. Linear models

k



L4.1 ANOVA



Theory









The Main Equation



$$d.f.(SST) = d.f.(SSTR) + d.f.(SSE)$$

 $n_T - 1 = (k - 1) + (n_T - k)$

Partitioning

The process of allocating the total sum of squares and degrees of freedom to the various components.







Example



$$SST = SSTR + SSE$$







Example

ANOVA table

A table used to summarize the analysis of variance computations and results. It contains columns showing the source of variation, the sum of squares, the degrees of freedom, the mean square, and the *F* value(s).

Let's perform for dataset 1: "good health"

depression2.txt

	Df	Sum Sq M	ean Sq	F value Pr(>F)	
Location	2	78.5	39.27	6.773 0.0023 **	
Residuals	57	330.4	5.80		
Signif. code	es:	0 `***′	0.001	`**' 0.01 `*' 0.05 `.' 0.1 ` ' 1	L







Some Definitions



Interaction

The effect produced when the levels of one factor interact with the levels of another factor in influencing the response variable.







ANOVA Table

Replications

The number of times each experimental condition is repeated in an experiment.

a = number of levels of factor A

- b = number of levels of factor B
- r = number of replications

 n_r = total number of observations taken in the experiment; $n_T = abr$

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Square	F
Factor A	SSA	<i>a</i> – 1	$MSA = \frac{SSA}{a-1}$	MSA MSE
Factor B	SSB	b - 1	$MSB = \frac{SSB}{b-1}$	$\frac{\text{MSB}}{\text{MSE}}$
Interaction	SSAB	(a - 1)(b - 1)	$MSAB = \frac{SSAB}{(a-1)(b-1)}$	MSAB MSE
Error	SSE	ab(r-1)	$MSE = \frac{SSE}{ab(r-1)}$	
Total	SST	$n_T - 1$		







Example



res = aov(Depression ~ Location + Health + Location*Health, Data)
summary(res)
source("http://sablab.net/scripts/drawANOVA.r")
drawANOVA(res)

ANOVA

ANOVA model:	D	epression = m + L	ocation + Health +	Location * Heal	th+e	
Factor	Df	Sum Sq	Mean Sq	F value	p-value	
Location	2	73.85	36.93	4.29	1.5981e-02	*
Health	1	1748.03	1748.03	203.094	4.3961e-27	***
Location:Health	2	26.12	13.06	1.517	2.2373e-01	











Experimental Design

Aware of Batch Effect !

When designing your experiment always remember about various factors which can effect your data: batch effect, personal effect, lab effect...











Experimental Design

Completely randomized design

An experimental design in which the treatments are randomly assigned to the experimental units.



We can nicely randomize:

Day effect

Batch effect







Experimental Design

Blocking

The process of using the same or similar experimental units for all treatments. The purpose of blocking is to remove a source of variation from the error term and hence provide a more powerful test for a difference in population or treatment means.









A good suggestion... ③

Block what you can block, randomize what you cannot, and try to avoid unnecessary factors









```
# L4.1 ANOVA
## clear memory
rm(list = ls())
## load data
Data = read.table("http://edu.sablab.net/data/txt/depression2.txt",header=T,sep="\t")
str(Data)
DataGH = Data[Data$Health == "good",]
## build 1-factor ANOVA model
res1 = aov(Depression ~ Location, DataGH)
summary(res1)
## build the ANOVA model
res2 = aov(Depression ~ Location + Health + Location*Health, Data)
## show the ANOVA table
summary(res2)
## Load function
source("http://sablab.net/scripts/drawANOVA.r")
x11()
drawANOVA(res2)
```





Dependent and Independent Variables







Dependent and Independent Variables







Example



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.

L4. Linear models





Regression Model and Regression Line

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 and Regression Line

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

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







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

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

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

Estimation

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 $y = b_0 + b_1 x$

cells.xls

1. Make a scatter plot for the data.



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







Overview







Slope and Intercept

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:

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

$$b_0 = m_y - b_1 m_x$$





Coefficient of Determination



The Main Equation







ANOVA and Regression



SST = SSTR + SSE



$$SST = SSR + SSE$$





Coefficient of Determination

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).











2014

32

Assumptions

Assumptions for Simple Linear Regression

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







Estimation of σ^2

SSE

i-th residual

The difference between the observed value of the dependent variable and the value predicted using the estimated regression equation; for the *i*-th observation the *i*-th residual is: $y_i - y_i$

Mean square error

The unbiased estimate of the variance of the error term σ^2 . It is denoted by MSE or s^2 . Standard error of the estimate: the square root of the mean square error, denoted by s. It is the estimate of σ , the standard deviation of the error term ε .







Sampling Distribution for b1

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

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

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





Test for Significance

- $H_0: \beta_1 = 0$ 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-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



1.	Calculate	manually	$y b_1$	and b_0
_				

Intercept	b0=	-191.008119
Slope	b1=	15.3385723

In Excel use the function:

= INTERCEPT(y, x)

2. Let's do it automatically

Data \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





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







Example2

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

Build the model and provide the p-value for linear dependency
 Provide interval estimation for the slope of the dependency
 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 $\beta_1 = 3$



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







Thank you for your attention

to be continued...



L4. Linear models

edu.sablab.net/abs2014 44