

# STATISTICAL DATA ANALYSIS IN EXCEL

# Lecture 6

# **Some Advanced Topics**

**Dr. Petr Nazarov** 

petr.nazarov@crp-sante.lu

31-10-2011

Statistical data analysis in Excel.

6. Some advanced topics



# **Correction for Multiple Comparisons**

Please download the data from edu.sablab.net/data/xls

all\_data.xls



#### **Correct Results and Errors**



Probability of an error in a multiple test:

1-(0.95)number of comparisons





#### False discovery rate (FDR)

CENTRE DE RECHERCHE

FDR control is a statistical method used in multiple hypothesis testing to correct for multiple comparisons. In a list of rejected hypotheses, FDR controls the expected proportion of incorrectly rejected null hypotheses (type I errors).

		Population Condition		
		H <sub>0</sub> is TRUE	H <sub>0</sub> is FALSE	Total
onclusion	Accept H <sub>0</sub> (non-significant)	U	T	m-R
	Reject H <sub>0</sub> (significant)	V	S	R
Ŭ	Total	$m_0$	$m-m_0$	т

$$FDR = E\left(\frac{V}{V+S}\right)$$



#### **False Discovery Rate**



Assume we need to perform k = 100 comparisons, and select maximum FDR =  $\alpha = 0.05$ 

#### Independent tests

The Simes procedure ensures that its expected value  $\mathbb{E}\left[\frac{V}{V+S}\right]$  is less than a given  $\alpha$  (Benjamini and Hochberg 1995). This procedure is valid when the *m* tests are independent. Let  $H_1 \dots H_m$  be the null hypotheses and  $P_1 \dots P_m$  their corresponding p-values. Order these values in increasing order and denote them by  $P_{(1)} \dots P_{(m)}$ . For a given  $\alpha$ , find the largest k such that  $P_{(k)} \leq \frac{k}{m} \alpha$ . Then reject (i.e. declare positive) all  $H_{(i)}$  for  $i = 1, \dots, k$ .

Note that the mean  $\alpha$  for these *m* tests is  $\frac{\alpha(m+1)}{2m}$  which could be used as a rough FDR, or RFDR, " $\alpha$  adjusted

for *m* indep. tests." The RFDR calculation shown here provides a useful approximation and is not part of the Benjamini and Hochberg method; see AFDR below.



#### **False Discovery Rate**

## Assume we need to perform k = 100 comparisons, and select maximum FDR = $\alpha = 0.05$

$$FDR = E\left(\frac{V}{V+S}\right)$$

Expected value for FDR <  $\alpha$  if

$$P_{(k)} \leq \frac{k}{m} \alpha$$



$$\frac{mP_{(k)}}{k} \leq \alpha$$



# **Example: Acute Lymphoblastic Leukemia**



Acute lymphoblastic leukemia (ALL), is a form of leukemia, or cancer of the white blood cells characterized by excess lymphoblasts.

**all\_data.xls** contains the results of full-trancript profiling for ALL patients and healthy donors using Affymetrix microarrays. The data were downloaded from ArrayExpress repository and normalized. The expression values in the table are in  $\log_2$  scale.

### Let us analyze these data:

- Calculate log-ratio (logFC) for each gene
- Calculate the p-value based on t-test for each gene
- Perform the FDR-based adjustment of the p-value.

Calculate the number of up and down regulated genes with FDR<0.01

How would you take into account logFC?

Example score:

 $score = -\log(adj.p.value) \cdot |logFC|$ 





look for "tetraspanin 7" + leukemia in google ③

Results are never perfect...



# **Empirical Interval Estimation for Random Functions**

Statistical data analysis in Excel. 6. Some advanced topics



## **Sum and Square of Normal Variables**

# Distribution of sum or difference of 2 normal random variables

The sum/difference of 2 (or more) normal random variables is a normal random variable with mean equal to sum/difference of the means and variance equal to SUM of the variances of the compounds.

$$x \pm y \rightarrow Normal \ distribution$$
$$E[x \pm y] = E[x] \pm E[y]$$
$$\sigma_{x \pm y}^{2} = \sigma_{x}^{2} + \sigma_{y}^{2}$$

**Distribution of sum of squares on** *k* **standard normal random variables** The sum of squares of *k* standard normal random variables is a  $\chi^2$  with *k* degree of freedom.

*if* 
$$x_1, ..., x_k \to Normal distribution$$
  
$$\sum_{i=1}^k x_i^2 \to \chi^2 \quad with \ d.f. = k$$

# What to do in more complex situations?

$$\frac{x}{y} \to ? \qquad \qquad \sqrt{x} \to ? \qquad \qquad \log(|x|) \to ?$$

Statistical data analysis in Excel.



# **Terrifying Theory**

# Try to solve analytically?

Simplest case. E[x] = E[y] = 0

## Ratio distribution

From Wikipedia, the free encyclopedia

A **ratio distribution** (or *quotient distribution*) is a probability distribution constructed as the distribution of the ratio of random variables having two other known distributions. Given two random variables X and Y, the distribution of the random variable Z that is formed as the ratio

$$Z = X/Y$$
is a ratio distribution. 
$$p_Z(z) = \frac{b(z) \cdot c(z)}{a^3(z)} \frac{1}{\sqrt{2\pi}\sigma_x \sigma_y} \left[ 2\Phi\left(\frac{b(z)}{a(z)}\right) - 1 \right] + \frac{1}{a^2(z) \cdot \pi\sigma_x \sigma_y} e^{-\frac{1}{2}\left(\frac{\mu_x^2}{\sigma_x^2} + \frac{\mu_y^2}{\sigma_y^2}\right)}$$

where

$$\begin{aligned} a(z) &= \sqrt{\frac{1}{\sigma_x^2} z^2 + \frac{1}{\sigma_y^2}} \\ b(z) &= \frac{\mu_x}{\sigma_x^2} z + \frac{\mu_y}{\sigma_y^2} \\ c(z) &= e^{\frac{1}{2} \frac{b^2(z)}{a^2(z)} - \frac{1}{2} \left(\frac{\mu_x^2}{\sigma_x^2} + \frac{\mu_y^2}{\sigma_y^2}\right)} \\ \Phi(z) &= \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}u^2} \ du \end{aligned}$$



## **Practical Approach**

Two rates where measured for a PCR experiment: experimental value (X) and control (Y). 5 replicates where performed for each.

From previous experience we know that the error between replicates is normally distributed.

**Q1:** provide an interval estimation for the fold change X/Y ( $\alpha$ =0.05)

**Q2:** provide an interval estimation for the log fold change  $log_2(X/Y)$ 

#	Experiment	Control
1	215	83
2	253	75
3	198	62
4	225	91
5	240	70
Mean	226.2	76.2
StDev	21.39	11.26

Let us use a *numerical simulation...* 



## **Practical Approach**

**1.** Generate 2 sets of 65536 normal random variable with means and standard deviations corresponding to ones of experimental and control set.

In Excel go: Tools  $\rightarrow$  Data Analysis:

Random Number Generation

If you do not have Data Analysis tool – approximate normal distribution by sum of uniform:

$$N(x,m_x,\sigma_x) = m_x + \sigma_x \left( \sum_{i=1}^{12} U(x_i) - 6 \right)$$

$$\bullet$$
 = RAND()  $\leftarrow$   $U(x)$ 

Mean	226.2	76.2
StDev	21.39	11.26

Number of Variables: 1	ОК
Number of Random Numbers: 65536	Cancel
Distribution: Normal	Help
Parameters	
M <u>e</u> an = 76.2	
Standard deviation = 11.26	
Random Seed:	
Output options	
Output Range: \$G:\$G     Signature	
New Worksheet Ply:	
O New <u>W</u> orkbook	



# **Practical Approach**

1. Generate 2 sets of 65536 normal random v with means and standard deviations correspo	Mean StDev	226.2 21.39	76.2 11.26	
ones of experimental and control set.	es of experimental and control set.		226.088799 21.379652	76.2823 11.2885
2. Build the target function. For Q1 build X/Y	X/Y.m X/Y.s min max	3.03289298 0.566865 -8.14098141 7.72162205		

**3.** Study the target function. Calculate summary, build histogram.



4. If you would like to have 95% interval, calculate 2.5% and 97.5% percentiles.
In Excel use function



**Practical Approach** 

What was a "mistake" in the previous case?



# There we spoke about prediction interval of X/Y. Now let's produce the interval estimation for mean X/Y

Mean	226.2	76.	2
StDev	9.57	5.0	3
X/Y.m	2.98047943		
X/Y.s	0.23616818		
min	2.01556098		
max	4.31131109		
		2 4 9 1	
ן <b>⊏נ∧/ ז</b>	$j \in [2.55,$	J.40 ]	





# **Practical Approach**

# **Q2:** provide an interval estimation for the log fold change log2(X/Y)



S	imulation	Normal
2.50%	1.3546	1.3482
97.50%	1.7998	1.7939



# Goodness of Fit and Independence

**Statistical data analysis in Excel.** 6. Some advanced topics



**Multinomial Population** 

#### **Multinomial population**

CENTRE DE RECHERCHE

A population in which each element is assigned to one and only one of several categories. The multinomial distribution extends the binomial distribution from two to three or more outcomes.

#### **Contingency table = Crosstabulation**

Contingency tables or crosstabulations are used to record, summarize and analyze the relationship between two or more categorical (usually) variables.

The new treatment for a disease is tested on 200 patien	ts.
The outcomes are classified as:	

- A patient is **completely treated**
- **B** disease transforms into a chronic form
- C treatment is unsuccessful 😕

In parallel the 100 patients treated with standard methods are observed

Category	Experimental	Control
A	94	38
В	42	28
С	64	34
Sum	200	100

The proportions for 3 "classes" of patients with and without treatment are:



Statistical data analysis in Excel.

6. Some adv



# **TEST OF GOODNESS OF FIT**

# **Goodness of Fit**

#### **Goodness of fit test**

A statistical test conducted to determine whether to reject a hypothesized probability distribution for a population.

**Model** – our assumption concerning the distribution, which we would like to test.

**Observed frequency** – frequency distribution for experimentally observed data,  $f_i$ 

**Expected frequency** – frequency distribution, which we would expect from our **model**,  $e_i$ 

#### Hypotheses for the test:

 $H_0$ : the population follows a multinomial distribution with the probabilities, specified by **model** 

 $H_a$ : the population does not follow ... model



Test statistics for goodness of fit



 $\chi^2$  has **k**-1 degree of freedom

At least 5 expected must be in each category!

#### Statistical data analysis in Excel. 6. Some advanced topics

# TEST OF GOODNESS OF FIT

# Example

The new treatment for a disease is tested on 200 patients. The outcomes are classified as:

- A patient is completely treated
- B disease transforms into a chronic form
- C treatment is unsuccessful (2)

In parallel the 100 patients treated with standard methods are observed

Model for

Expected

#### **1.** Select the model and calculate expected frequencies

Let's use control group (classical treatment) as a model, then:

Control

Category

5	frequencies	control	freq., e	
А	38	0.38	76	
В	28	0.28	56	
С	34	0.34	68	
Sum	100	1	200	

CHIDIST( $\chi^2$ ,d.f.)

$$\bullet$$
 = CHITEST( $f, e$ )

3. Calculate

d.f. = *k*–1

p-value for  $\chi^2$  with

**2.** Compare expected frequencies with the experimental ones and build  $\chi^2$ 

 $\chi^2 = \sum^k \frac{(f_i - e_i)^2}{(f_i - e_i)^2}$ 

(f-e)2/e

4.263

3.500

0.235

7.998

Experimental

freq., f

94

42

64

200

Category

Α

В

С

Chi<sub>2</sub>











#### **Goodness of Fit for Independence Test: Example**

Alber's Brewery manufactures and distributes three types of beer: white, regular, and dark. In an analysis of the market segments for the three beers, the firm's market research group raised the question of whether preferences for the three beers differ among male and female beer drinkers. If beer preference is independent of the gender of the beer drinker, one advertising campaign will be initiated for all of Alber's beers. However, if beer preference depends on the gender of the beer drinker, the firm will tailor its promotions to different target markets.

#### beer.xls





 $H_a$ : Beer preference is **not independent** of the gender of the beer drinker

sex\beer	White	Regular	Dark	Total
Male	20	40	20	80
Female	30	30	10	70
Total	50	70	30	150







# **Goodness of Fit for Independence Test: Example**

1 Ruild model	sex\beer	White	Regular	Dark	Total
	Male	20	40	20	80
assuming	Female	30	30	10	70
independence	Total	50	70	30	150

	White	Regular	Dark	Total
Model	0.3333	0.4667	0.2000	1

2. Transfer the model into expected frequencies, multiplying model value by number in group

sex\beer	White	Regular	Dark	Total
Male	26.67	37.33	16.00	80
Female	23.33	32.67	14.00	70
Total	50	70	30	150

$$e_{ij} = \frac{(Row \ i \ Total)(Column \ j \ Total)}{Sample \ Size}$$

#### **3.** Build $\chi^2$ statistics



 $\chi^2$  distribution with d.f.=(n-1)(m-1), provided that the expected frequencies are 5 or more for all categories.

#### 4. Calculate p-value

p-value = 0.047, reject H<sub>0</sub>

Statistical data analysis in Excel.



## **Test for Normality: Example**

Chemline hires approximately 400 new employees annually for its four plants. The personnel director asks whether a normal distribution applies for the population of aptitude test scores. If such a distribution can be used, the distribution would be helpful in evaluating specific test scores; that is, scores in the upper 20%, lower 40%, and so on, could be identified quickly. Hence, we want to test the null hypothesis that the population of test scores has a normal distribution. The study will be based on 50 results.



Statistical data analysis in Excel.



# **TEST FOR CONTINUOUS DISTRIBUTIONS**

# **Test for Normality: Example**







# Thank you for your attention

