

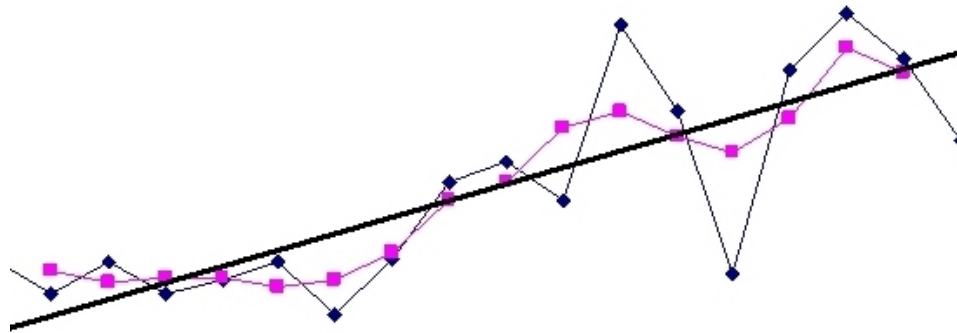
Mathematical and Statistical Modelling in Medicine

Author: Tibor Nyári PhD

University of Szeged
Department of Medical Physics and Informatics

www.model.u-szeged.hu
www.szote.u-szeged.hu/dmi

Analysis of variance



Method

- If the null hypothesis is true, then the populations are the same: they are normal, and they have the same mean and the same variance. We will estimate the numerical value of this common variance in two distinct ways: we will compute the “between-groups variance” and the “within-groups variance”.
- If the null hypothesis is true, then these two distinct estimates of the variance should be equal, their equality can be tested by an F ratio test. The fact that we can compute variances in two ways stems from the break-down of the total sum of squares into a “between-groups sum of squares” and a “within-groups sum of squares”.
- The total number of degrees of freedom $N-1$, where N is the sum of all sample sizes, is also broken down into the appropriate “between-groups” and “within-groups” degrees of freedom: $t-1$ and $N-t$, respectively. The results of computations of ANOVA methods are usually tabulated. The rows of such tables give the source of the variance; the columns contain the sum of squares, the number of degrees of freedom, the variances, the F -value (variance ratio), and the p -value.

ANOVA table

Source of variation	Sum of squares	Degrees of freedom	Variance	F
Between groups	$Q_b = \sum_{i=1}^h n_i (\bar{x}_i - \bar{\bar{x}})^2$	h-1	$s_b^2 = \frac{Q_b}{h-1}$	$F = \frac{s_b^2}{s_w^2}$
Within groups	$Q_w = \sum_{i=1}^h \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2$	N-h	$s_w^2 = \frac{Q_w}{N-h}$	
Total	$Q = \sum_{i=1}^h \sum_{j=1}^{n_i} (x_{ij} - \bar{\bar{x}})^2$	N-1		

Multiple comparisons

- If the result of the ANOVA is not significant at the specified level, the analysis is complete. We expect all samples be drawn from the same population; the differences between sample means are due to random effects
- If the result of the ANOVA is significant, then we have to accept the alternative hypothesis: there is at least one group different from one of the others.
- To find these groups, we have to compare each group with each of the others. As the two-sample t -test is inappropriate to do this, there are special tests for multiple comparisons that keep the probability of Type I error as α .

Pairwise methods I.

- The most often used multiple comparisons are the modified t-tests. Another often used method is the so-called Bonferroni method: to achieve a level of not more than α for a set of a number of c tests, we need to choose a level α / c for the individual tests.

$$t = \frac{\bar{x}_i - \bar{x}_j}{s_b \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}}$$

- For example for three comparisons a p-value less than $0.05/3=0.017$ has to be considered significant instead of $p=0.05$. This method is conservative. We know only that the probability does not exceed α for the set. This method can be used in cases involving small numbers of comparisons.

Scheffé test

- The *Scheffé* test performs simultaneous joint pairwise comparisons for all possible pairwise combinations of means. This test can be used to examine special linear combinations of group means, not simply pairwise comparisons.

$$\left(\frac{\bar{x}_i - \bar{x}_j}{s_b \sqrt{\left(\frac{1}{n_i} + \frac{1}{n_j}\right)(h - 1)}} \right)^2$$

Tukey-test (only in the case of equal sample sizes)

- The following quantity has to be compared with a value from the special table of studentized ranges

$$T = \frac{\overline{x_i} - \overline{x_j}}{\frac{s_b}{\sqrt{n}}}$$

Dunnett test

- Another multiple comparison method is the *Dunnett test*: a test comparing a given group (control) with the others.

Example

- Assume that we have recorded the biomass of 3 bacteria in flasks of glucose broth, and we used 3 replicate flasks for each bacterium.

Replicate	Bacterium A	Bacterium B	Bacterium C
1	12	20	40
2	15	19	35
3	9	23	42

Anova: Single Factor

Descriptive statistics

<i>Groups</i>	<i>Count</i>	<i>Sum</i>	<i>Average</i>	<i>Variance</i>
Bacterium A	3	36	12	9
Bacterium B	3	62	20,66667	4,333333
Bacterium C	3	117	39	13

ANOVA table

<i>Source of Variation</i>	<i>SS</i>	<i>df</i>	<i>MS</i>	<i>F</i>	<i>P-value</i>	<i>F crit</i>
Between Groups	1140,222	2	570,1111	64,94937	8,61E-05	5,143253
Within Groups	52,66667	6	8,777778			
Total	1192,889	8				

Two-way analysis of variance

Main Effect

- The main effect involves the independent variables one at a time. The interaction is ignored for this part. Just the rows or just the columns are used, not mixed. This is the part which is similar to the one-way analysis of variance. Each of the variances calculated to analyze the main effects are like the between variances

Interaction Effect

- The interaction effect is the effect that one factor has on the other factor. The degrees of freedom here is the product of the two degrees of freedom for each factor.

Within Variation

- The Within variation is the sum of squares within each treatment group. You have one less than the sample size (remember all treatment groups must have the same sample size for a two-way ANOVA) for each treatment group. The total number of treatment groups is the product of the number of levels for each factor. The within variance is the within variation divided by its degrees of freedom. The within group is also called the error.

F-Tests

- There is an F-test for each of the hypotheses, and the F-test is the mean square for each main effect and the interaction effect divided by the within variance. The numerator degrees of freedom come from each effect, and the denominator degrees of freedom is the degrees of freedom for the within variance in each case.

Hypothesis

- There are three sets of hypothesis with the two-way ANOVA.
- The null hypotheses for each of the sets are given below.
- The population means of the first factor are equal. This is like the one-way ANOVA for the row factor.
- The population means of the second factor are equal. This is like the one-way ANOVA for the column factor.
- There is no interaction between the two factors. This is similar to performing a test for independence with contingency tables.

ANOVA table without replication

$$T_j = \sum_{i=1}^p x_{ij}$$

$$Z_i = \sum_{j=1}^q x_{ij}$$

$$T = \sum_{j=1}^q T_j$$

$$\bar{x}_{i\bullet} = \frac{Z_i}{q}$$

$$\bar{x}_{\bullet j} = \frac{T_j}{p}$$

		Factor B				Z_i	$\bar{x}_{i\bullet}$
		j=1	j=2	j=q		
Factor A	i=1	x_{11}	x_{12}	x_{1q}	Z_1	$\bar{x}_{1\bullet}$
	i=2	x_{21}	x_{22}	x_{2q}	Z_2	$\bar{x}_{2\bullet}$
	
	
	
	i=p	x_{k1}	x_{k2}	x_{pq}	Z_p	$\bar{x}_{p\bullet}$
	T_j	T_1	T_2	T_q	T	

$\bar{x}_{\bullet j}$	$\bar{x}_{\bullet 1}$	$\bar{x}_{\bullet 2}$		$\bar{x}_{\bullet q}$	$\bar{\bar{x}}$
-----------------------	-----------------------	-----------------------	--	-----------------------	-----------------

$$\text{Intercept / correction factor} = \frac{T^2}{pq}$$

ANOVA table without replication

Source of variation	Sum of squares	Degrees of freedom	Variance (Var=Q/f)	F - value (F=Var/Var _{error})
Correction factor	Q_{CF}	$f_{CF} = 1$	$Var_{CF} = Q_{CF}/f_{CF}$	$F_{CF} = Var_{CF}/Var_{error}$
Faktor A	Q_A	$f_A = p - 1$	$Var_A = Q_A/f_A$	$F_A = Var_A/Var_{error}$
Faktor B	Q_B	$f_B = q - 1$	$Var_B = Q_B/f_B$	$F_B = Var_B/Var_{error}$
Error	Q_{error}	$f_{error} = (p-1)(q-1)$	$Var_{error} = Q_{error}/f_{error}$	
Total	Q_{total}	$f_{total} = pq - 1$		

ANOVA table without replication

Source of variation	Sum of squares (SS)	Degrees of freedom (df)	Variance (MS=SS/df)	F
Intercept	$SS_{CF} = \frac{T^2}{pq}$	1	S_{CF}^2	$F = \frac{S_{CF}^2}{S_h^2}$
Factor A	$SS_A = \left(\frac{1}{q} \sum_i Z_i^2 \right) - \left(\frac{T^2}{pq} \right)$	p-1	S_A^2	$F = \frac{S_A^2}{S_h^2}$
Factor B	$SS_B = \left(\frac{1}{p} \sum_j T_j^2 \right) - \left(\frac{T^2}{pq} \right)$	q-1	S_B^2	$F = \frac{S_B^2}{S_h^2}$
Error (Residual)	$SS_{Error} = SS_{Total} - SS_A - SS_B$	(p-1)(q-1)	S_h^2	
Total	$SS_{Total} = \left(\sum_{i,j} x_{ij}^2 \right) - \left(\frac{T^2}{pq} \right)$	N-1		

ANOVA table with replication

Source of variation	Sum of squares	Degrees of freedom	Variance (Var=Q/f)	F - value (F=Var/Var _{error})
Correction factor	Q_{CF}	$f_{CF} = 1$	$Var_{CF} = Q_{CF}/f_{CF}$	$F_{CF} = Var_{CF}/Var_{error}$
Faktor A	Q_A	$f_A = p - 1$	$Var_A = Q_A/f_A$	$F_A = Var_A/Var_{error}$
Faktor B	Q_B	$f_B = q - 1$	$Var_B = Q_B/f_B$	$F_B = Var_B/Var_{error}$
A x B	Q_{AB}	$f_{AB} = (p-1)(q-1)$	$Var_{AB} = Q_{AB}/f_{AB}$	$F_{AB} = Var_{AB}/Var_{error}$
Error	Q_{error}	$f_{error} = N - pq (=0)$	$Var_{error} = Q_{error}/f_{error}$	
Total	Q_{total}	$f_{total} = pq - 1$		

ANOVA table with replication

Source of variation	Sum of squares (SS)	Degrees of freedom (df)	Variance (MS=SS/df)	F
Intercept	$SS_{CF} = \frac{T^2}{pq}$	1	S_{CF}^2	$F = \frac{S_{CF}^2}{S_h^2}$
Factor A	$SS_A = nq \sum_{i=1}^p (\bar{x}_{i.} - \bar{x})^2$	p-1	S_A^2	$F = \frac{S_A^2}{S_h^2}$
Factor B	$SS_B = np \sum_{j=1}^q (\bar{x}_{.j} - \bar{x})^2$	q-1	S_B^2	$F = \frac{S_B^2}{S_h^2}$
Interaction (AxB)	$SS_{AxB} = n \sum_{i=1}^p \sum_{j=1}^q (\bar{x}_{ij} - \bar{x}_{i.} - \bar{x}_{.j} + \bar{x})^2$	(p-1)(q-1)	S_{AB}^2	$F = \frac{S_{AB}^2}{S_h^2}$
Error (Within)	$SS_{within} = \sum_{h=1}^n \sum_{i=1}^p \sum_{j=1}^q (x_{hij} - \bar{x}_{ij})^2$	N-pq	S_h^2	
Total	$SS_{total} = \sum_{h=1}^n \sum_{i=1}^p \sum_{j=1}^q (x_{hij} - \bar{x})^2$	N-1		

The two-way ANOVA model

- Let us denote the numbers of levels of factors 1 and 2 by t and l , respectively, and by N the total number of observations. The two-way ANOVA model is

$$y_{ijk} = \mu + \alpha_i + \beta_j + \theta_{ij} + \varepsilon_{ijk}, \quad i = 1, \dots, t \quad j = 1, \dots, l, \quad k = 1, \dots, n_{ij}$$

where we use the following notations:

- y_{ijk} = the k -th observed value of the dependent variable when we are using level i of factor 1 and level j of factor 2,
- μ = an overall mean, (unknown constant)
- α_i = the effect due to level i of factor 1 (an unknown constant),
- β_j = the effect due to level j of factor 2, (an unknown constant),
- θ_{ij} = the effect due to the interaction of level i of factor 1 and level j of factor 2 (an unknown constant),
- ε_{ijk} = the k -th error term when we are using level i of factor 1 and level j of factor 2 (assumed to be distributed as $N(0, \sigma)$)

According to the above questions, the following null hypotheses can be tested

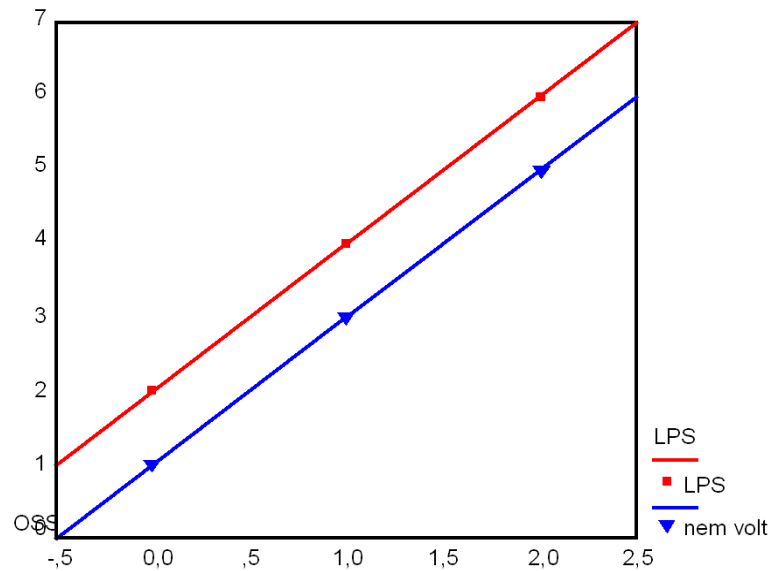
- $H_1: \alpha_1 = \alpha_2 = \dots = \alpha_t$
- $H_2: \beta_1 = \beta_2 = \dots = \beta_t$
- $H_3: \theta_{ij} = 0$

Example:

- Suppose that factor A has 3 levels ($i=3$) and factor B has 2 levels ($j=2$).
- Then theoretically, we might have the following situations

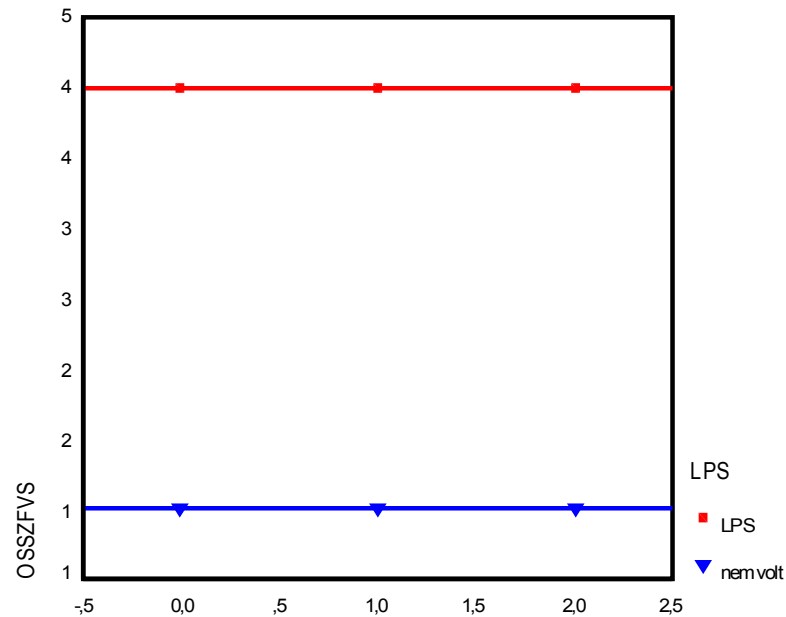
Only the effect of A is significant, the only significant factor is A

$$y_{ijk} = \mu + \alpha_i + \varepsilon_{ijk},$$



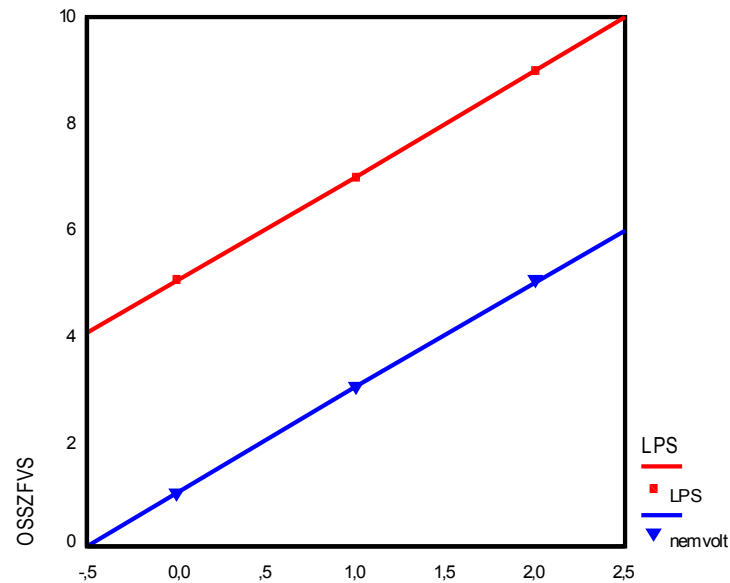
Only the effect of B is significant, the only significant factor is B

$$y_{ijk} = \mu + \beta_j + \varepsilon_{ijk}$$

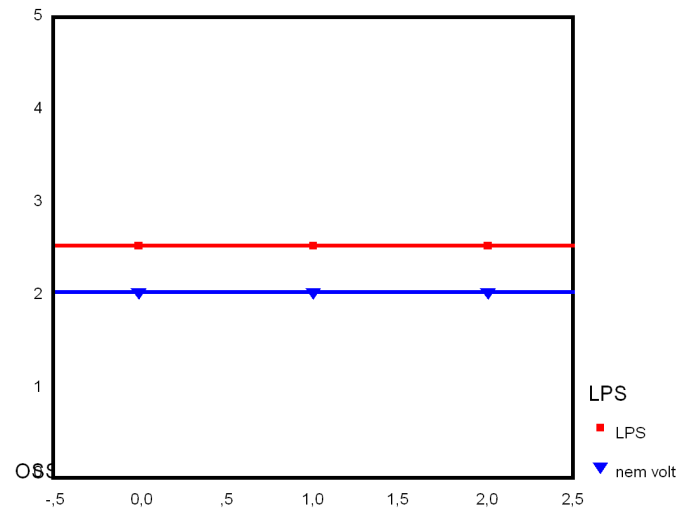


Both factors are significant

$$y_{ijk} = \mu + \alpha_i + \beta_j + \varepsilon_{ijk}$$



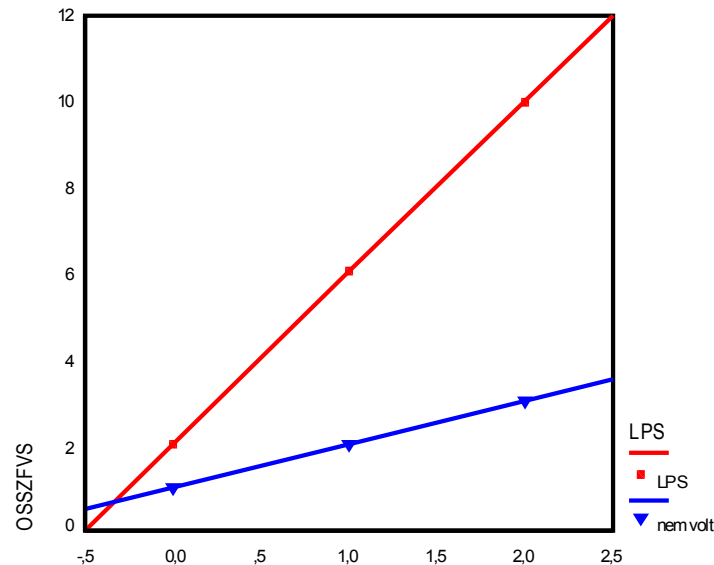
None of them is significant



Interaction

- The differences in A depend on the level of factor B

$$y_{ijk} = \mu + \alpha_i + \beta_j + \Theta_{ij} + \varepsilon_{ijk}, \quad i = 1, \dots, t \quad j = 1, \dots, l,$$



Two-way ANOVA

- In two-way ANOVA, the total sum of squares is decomposed into four terms, according to the effects in the model.
- The results are generally written into an ANOVA table which contains rows for the effects of factors 1 and 2, the interaction and the error with $(p-1)$, $(q-1)$, $(p-1)(q-1)$ and $(N-pq)$ degrees of freedom, respectively

- The rows of this tables give the components for the effects of factor 1, factor 2, the interaction and the error term, while the columns contain the sum of squares, the number of degrees of freedom, the variances, the F -values (variance ratio), and the p -value of F .

Decision

- There are three F -values in this table according to the three hypotheses.
- The interaction, is tested first. If it is not significant, the significance of each of factors 1 and 2 can be tested separately.
- If H_1 is rejected, we can say that at least two means of the factor 1 differ.
- If t (which is more than two), the number of levels of factor 1, we again have to use multiple comparisons to find pairwise differences.

Decision

- If the interaction is significant, then the relationship between the means of factor 1 depends on the level of factor 2. Multiple comparisons can be performed for each combination of one factor with a given level of the other factor.
- There are special methods against the increase of Type I error, and the use of t -tests independently is an incorrect solution

Example

- Suppose that we have grown one bacterium in broth culture at 3 different pH levels at 4 different temperatures. We have 12 flasks in all, but no replicates. Growth was measured by optical density (O.D.).

Anova: Two-Factor Without Replication DATA

Temp °C	pH 5.5	pH 6.5	pH 7.5
25	10	19	40
30	15	25	45
35	20	30	55
40	15	22	40

Descriptive statistics

<i>SUMMARY</i>	<i>Count</i>	<i>Sum</i>	<i>Average</i>	<i>Variance</i>
25	3	69	23	237
30	3	85	28,33333	233,3333
35	3	105	35	325
40	3	77	25,66667	166,3333
pH 5.5	4	60	15	16,66667
pH 6.5	4	96	24	22
pH 7.5	4	180	45	50

Two-way Anova

<i>Source of Variation</i>	<i>SS</i>	<i>df</i>	<i>MS</i>	<i>F</i>	<i>P-value</i>	<i>F crit</i>
Rows	238,667	3	79,5555	17,4634	0,00228	4,75706
Columns	1896	2	948	208,19	2,8E-06	5,143
Error	27,3333	6	4,55556			
Total	2162	11				

- Of interest, another piece of information is revealed by this analysis - the effects of temperature do not interact with effects of pH. In other words, a change of temperature does not change the response to pH, and vice-versa. We can deduce this because the residual (error) mean square (MS) is small compared with the mean squares for temperature (columns) or pH (rows). [A low residual mean square tells us that most variation in the data is accounted for by the separate effects of temperature and pH].

Example II

Temp °C	pH 5.5	pH 6.5	pH 7.5
25	20	38	80
30	30	50	60
35	40	60	50
40	50	44	20

ANOVA without replication

<i>Source of Variation</i>	<i>SS</i>	<i>df</i>	<i>MS</i>	<i>F</i>	<i>P-value</i>	<i>F crit</i>
Rows	233	3	77,66667	0,193682	0,897001	4,757063
Columns	660,6667	2	330,3333	0,823774	0,482933	5,143253
Error	2406	6	401			
Total	3299,667	11				

ANOVA with replication

- The data of serum IGG level were measured in triplicate of three patients.
- Is there any difference in individual mean IGG level values using different method of measurement?

ANOVA with replication Dataset

help	species	method	value
a1	A	1	23.8
a2	A	1	11
a3	A	1	25.2
b1	B	1	9.7
b2	B	1	4.8
b3	B	1	5.2
c1	C	1	15.2
c2	C	1	10.4
c3	C	1	12.8
a1	A	2	10.7
a2	A	2	8.6
a3	A	2	10.6
b1	B	2	1.9
b2	B	2	3.2
b3	B	2	1.2
c1	C	2	9.7
c2	C	2	7.8
c3	C	2	6.5

The role of Column „help”

- Using Excel Pivot Table command to get the „ANOVA” table we need to use an extra column let’s name it as „help”.
- Without this we get the following table (there is no replicant).

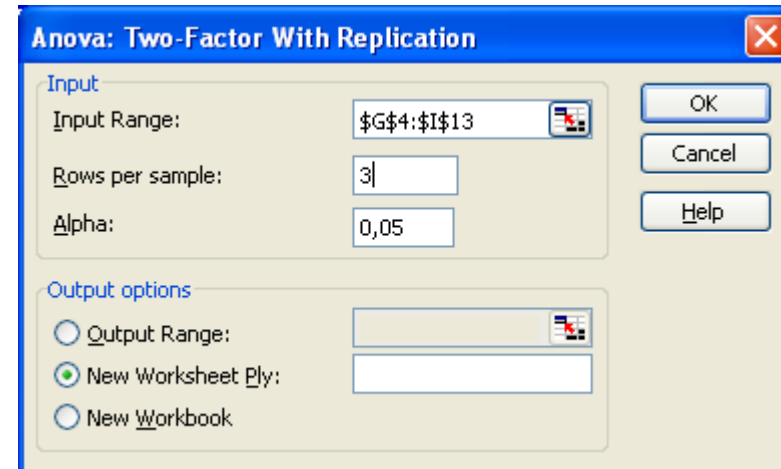
Sum of value	method		
species	1	2	Grand Total
A	60	29.9	89.9
B	19.7	6.3	26
C	38.4	24	62.4
Grand Total	118.1	60.2	178.3

Pivot table using „help” column

Sum of value	method		
help	1	2	Grand Total
a1	23.8	10.7	34.5
a2	11	8.6	19.6
a3	25.2	10.6	35.8
b1	9.7	1.9	11.6
b2	4.8	3.2	8
b3	5.2	1.2	6.4
c1	15.2	9.7	24.9
c2	10.4	7.8	18.2
c3	12.8	6.5	19.3
Grand Total	118.1	60.2	178.3

The settings of ANOVA command

	G	H	I
4	Species	Method_1	Method_2
5	A	23.8	10.7
6	rep2	11	8.6
7	rep3	25.2	10.6
8	B	9.7	1.9
9	rep2	4.8	3.2
10	rep3	5.2	1.2
11	C	15.2	9.7
12	rep2	10.4	7.8
13	rep3	12.8	6.5



Anova: Two-Factor With Replication I.

SUMMARY	Method_1	Method_2	Total
<i>A</i>			
Count	3	3	6
Sum	60	29.9	89.9
Average	20	9.966667	14.98333
Variance	61.24	1.403333	55.25767
<i>B</i>			
Count	3	3	6
Sum	19.7	6.3	26
Average	6.566667	2.1	4.333333
Variance	7.403333	1.03	9.358667
<i>C</i>			
Count	3	3	6
Sum	38.4	24	62.4
Average	12.8	8	10.4
Variance	5.76	2.59	10.252

Anova: Two-Factor With Replication II.

ANOVA

<i>Source of Variation</i>	<i>SS</i>	<i>df</i>	<i>MS</i>	<i>F</i>	<i>P-value</i>	<i>F crit</i>
Sample/Species	342.4678	2	171.2339	12.93524	0.001012	3.885294
Columns/Meth	186.245	1	186.245	14.0692	0.002766	4.747225
Interaction	29.24333	2	14.62167	1.104541	0.36282	3.885294
Within	158.8533	12	13.23778			
Total	716.8094	17				

SPSS results I

Descriptive Statistics

Dependent Variable: value

species	method	Mean	Std. Deviation	N
A	Method1	20,0000	7,82560	3
	Method2	9,9667	1,18462	3
	Total	14,9833	7,43355	6
B	Method1	6,5667	2,72091	3
	Method2	2,1000	1,01489	3
	Total	4,3333	3,05919	6
C	Method1	12,8000	2,40000	3
	Method2	8,0000	1,60935	3
	Total	10,4000	3,20187	6
Total	Method1	13,1222	7,24530	9
	Method2	6,6889	3,71835	9
	Total	9,9056	6,49348	18

Levene's Test of Equality of Error Variances^a

Dependent Variable: value

F	df1	df2	Sig.
6,402	5	12	,004

Tests the null hypothesis that the error variance of the dependent variable is equal across groups.

a. Design: Intercept+species+method+species * method

SPSS results II.

Tests of Between-Subjects Effects

Dependent Variable: value

Source	Type III Sum of Squares	df	Mean Square	F	Sig.
Corrected Model	557,956 ^a	5	111,591	8,430	,001
Intercept	1766,161	1	1766,161	133,418	,000
species	342,468	2	171,234	12,935	,001
method	186,245	1	186,245	14,069	,003
species * method	29,243	2	14,622	1,105	,363
Error	158,853	12	13,238		
Total	2482,970	18			
Corrected Total	716,809	17			

a. R Squared = ,778 (Adjusted R Squared = ,686)

Post-Hoc test

Multiple Comparisons

Dependent Variable: value

LSD

(I) species	(J) species	Mean Difference (I-J)	Std. Error	Sig.	95% Confidence Interval	
					Lower Bound	Upper Bound
A	B	10,6500*	2,10062	,000	6,0731	15,2269
	C	4,5833*	2,10062	,050	,0065	9,1602
B	A	-10,6500*	2,10062	,000	-15,2269	-6,0731
	C	-6,0667*	2,10062	,014	-10,6435	-1,4898
C	A	-4,5833*	2,10062	,050	-9,1602	-,0065
	B	6,0667*	2,10062	,014	1,4898	10,6435

Based on observed means.

*. The mean difference is significant at the ,05 level.