

# Chapter 9: Analysis of Variance (ANOVA)

## Contents

<b>1</b>	<b>One-way ANOVA – Introductory example</b>	<b>2</b>
1.1	Descriptive statistics . . . . .	2
1.2	Problem formulation . . . . .	3
<b>2</b>	<b>F-test for multiple means in one-way ANOVA</b>	<b>4</b>
2.1	ANOVA – testing principle . . . . .	4
2.2	Example . . . . .	6
<b>3</b>	<b>One-way ANOVA as a linear model</b>	<b>7</b>
3.1	Use of treatment coding . . . . .	7
3.2	Use of sum coding . . . . .	8
<b>4</b>	<b>Model diagnostics</b>	<b>10</b>
4.1	Assumptions made . . . . .	10
4.2	Checking assumptions . . . . .	10
4.3	Check assumption of homogeneity of variance . . . . .	10
4.4	Check assumption of normality . . . . .	11
4.5	Checking for influential observations . . . . .	13
<b>5</b>	<b>Pairwise comparisons of treatment effects</b>	<b>14</b>
5.1	Questions post-hoc . . . . .	14
5.2	Planned comparison of means in ANOVA . . . . .	15
5.3	Multiple comparisons . . . . .	16
5.3.1	The problem of multiple comparisons . . . . .	16
5.3.2	An overview of multiple comparisons procedures . . . . .	16
<b>6</b>	<b>Extensions</b>	<b>21</b>
6.1	The Kruskal-Wallis test . . . . .	21
<b>7</b>	<b>Two-way ANOVA</b>	<b>23</b>
7.1	Introduction . . . . .	23
7.2	Two-way ANOVA model . . . . .	27
7.3	Strategy for the analysis of two-way ANOVA studies . . . . .	28
7.3.1	Example in R . . . . .	28
7.4	Diagnostics . . . . .	29
7.4.1	Checking homogeneity of variances . . . . .	29
7.4.2	Checking normality of residuals . . . . .	30
7.4.3	Influential observations . . . . .	31
7.5	Multiple comparisons for the main effects (in case interaction is not significant) . . . . .	33
7.6	Two-way ANOVA when cells have unequal sample size . . . . .	34
7.6.1	What is an unbalanced design? . . . . .	34
7.6.2	Illustrative example . . . . .	34
7.6.3	ANOVA table . . . . .	36
7.6.4	Diagnostics . . . . .	37
7.6.5	Pairwise comparisons of treatment effects . . . . .	39

<b>8</b>	<b>Experimental design</b>	<b>41</b>
8.1	Observational study versus designed experiment . . . . .	41
8.2	Basic principles of experimental design . . . . .	42
8.2.1	Replication . . . . .	42
8.2.2	Randomization . . . . .	43
8.2.3	Blocking . . . . .	43
<b>9</b>	<b>The general linear model</b>	<b>44</b>

## 1 One-way ANOVA – Introductory example

A t-test compares means of two independent groups ( $H_0 : \mu_1 = \mu_2$ ). **One-way analysis of variance (one-way ANOVA) is a testing procedure to compare the means of multiple groups.**

### Example *Pollution*

In some cities in the USA, they collect measurements about pollution. Next variables are available:

Variable	Description
city	Name of the city (located in USA)
regio	Region in USA: W: West N: North NO: North-East ZO: South-East C: Central
JanT	Average temperature in January (in Fahrenheit)
JulT	Average temperature in July (in Fahrenheit)
Hum	Relative humidity (in percentages)
Rain	Yearly amount of rain (in inches)
Mortality	Mortality, corrected for age
Educ	Median of the education level
density	Density
NW	Percentage of non-white

Import the data set *pollutie2.txt* as *pollution*.

```
pollution <- read.table(file=file.choose(), header=TRUE)
summary(pollution)
head(pollution, n = 15)
```

We will consider the pollution example and look at the average yearly amount of rain across different regions.

### 1.1 Descriptive statistics

```
install.packages("psych")
library(psych)

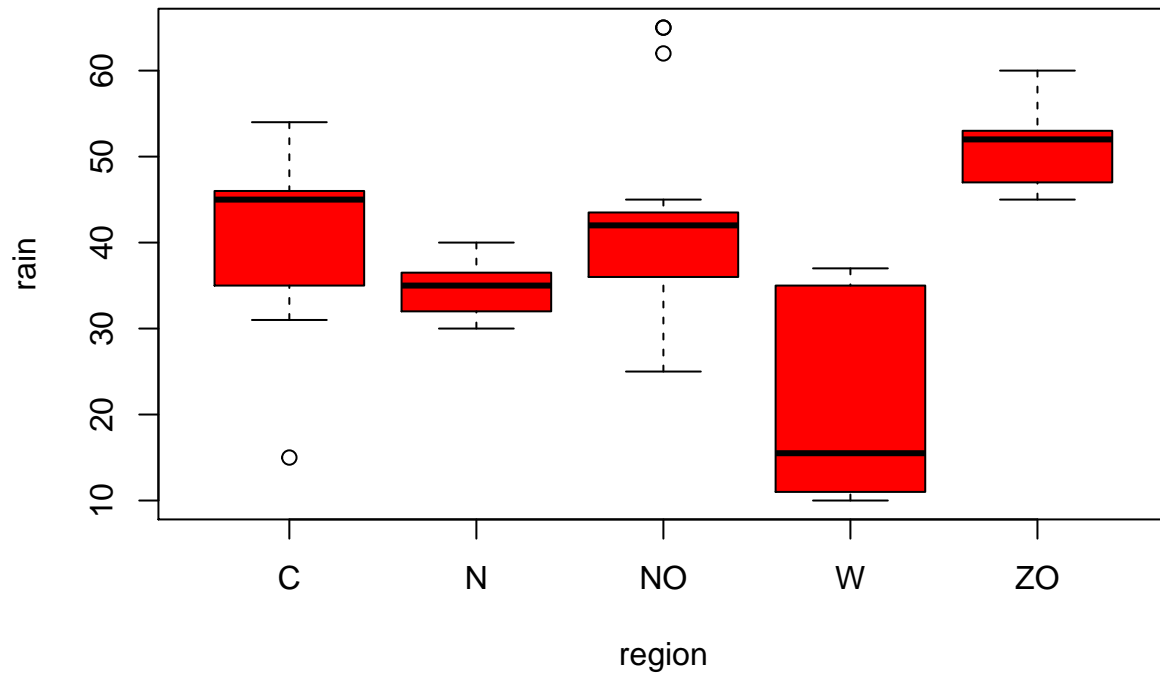
rain <- pollution$Rain
region <- pollution$regio

describe <- describeBy(rain, region, mat = TRUE)
describe.st <- subset(describe, select = c("group1", "n", "mean", "sd", "median", "min", "max"))
describe.st
```

##	group1	n	mean	sd	median	min	max
##	X11	C	9	40.55556	11.886033	45.0	15 54
##	X12	N	15	34.80000	3.098387	35.0	30 40
##	X13	NO	24	41.95833	9.993385	42.0	25 65

```
## X14      W 6 20.66667 12.209286 15.5 10 37
## X15      ZO 5 51.40000 5.856620 52.0 45 60
```

```
boxplot(rain ~ region, col = 2, names = levels(pollution$region))
```



### Remark:

Instead of using the `DescribeBy` function in the `psych` package, we can also use the `summarise` and `group_by` function from the `tidyverse` package.

Another way to obtain descriptive statistics:

```
library(tidyverse)
by_region <- group_by(pollution, regio)
summarise(by_region, Avgrain = mean(Rain))
```

```
## # A tibble: 5 x 2
##   regio Avgrain
##   <fct>   <dbl>
## 1 C      40.6
## 2 N      34.8
## 3 NO     42.0
## 4 W      20.7
## 5 ZO     51.4
```

## 1.2 Problem formulation

From the descriptive statistics analysis, we calculated the mean amount of rain for samples drawn from 5 different regions.

Region	Sample mean amount of rain	Population mean amount of rain
C	40.56	$\mu_1$
N	34.80	$\mu_2$
NO	41.96	$\mu_3$
W	20.67	$\mu_4$
ZO	51.40	$\mu_5$

We want to evaluate the following questions:

1. Do the five regions have the same average amount of rain ( $H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5$ )?
2. Which groups of regions have the same average amount of rain (homogeneous groups)?

Since there are more than two independent groups, a t-test cannot be used to compare the means. Instead, one-way ANOVA will be used to assess these questions.

## 2 F-test for multiple means in one-way ANOVA

### 2.1 ANOVA – testing principle

The one-way ANOVA compares the means between different groups. If there are  $r$  number of groups, then the ANOVA tests the following **hypothesis** :

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_r$$

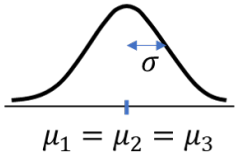
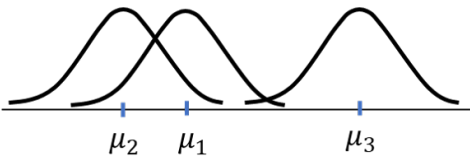
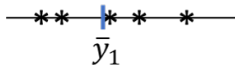
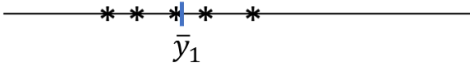
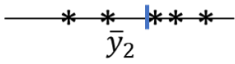
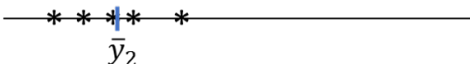
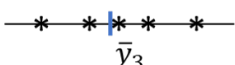
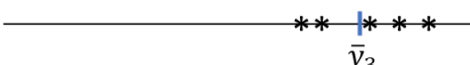
versus

$H_1$  : the means are not all the same

with  $\mu_i$  the population group mean of group  $i$ .

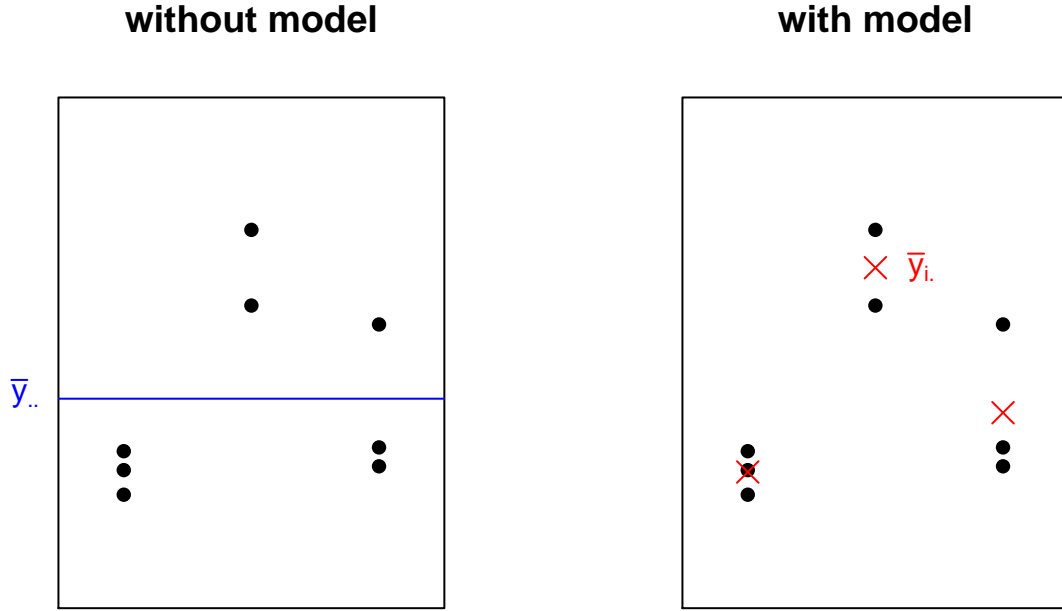
Consider a situation where the mean of three different groups are compared. ANOVA is used to test the null hypothesis  $H_0 : \mu_1 = \mu_2 = \mu_3$  against the alternative hypothesis that the means are not all the same.

We assume that the three groups correspond to 3 normally distributed populations with same variance  $\sigma^2$

	$H_0$ true	$H_0$ not true
<b>Population</b>		
<b>Sample</b>		
From population 1:		
From population 2:		
From population 3:		

**Testing principle:** Reject  $H_0$  if the variability of  $\bar{y}_i$  is too big compared to the within-group variance  $\sigma^2$ .

*Partitioning of the variances*



Total deviation:  $y_{ij} - \bar{y}_{..}$

Residual deviation:  $y_{ij} - \bar{y}_i$

Partitioning for observation  $ij$ :

$$(y_{ij} - \bar{y}_{..}) = (y_{ij} - \bar{y}_{i\bullet}) + (\bar{y}_{i\bullet} - \bar{y}_{..}) \quad (1)$$

(deviation of observation  $ij$ ) = (random error) + (effect from group  $i$ )

with

- $y_{ij}$  the value of the observation  $ij$
- $\bar{y}_{..}$  the mean of all the observations of all the groups
- $\bar{y}_{i\bullet}$  the mean of all the observation in group  $i$

**Sum of squares (SS) decomposition:**

$$\sum_{i=1}^r \sum_{j=1}^{n_j} (y_{ij} - \bar{y}_{..})^2 = \sum_{i=1}^r \sum_{j=1}^{n_j} (y_{ij} - \bar{y}_{i\bullet})^2 + \sum_{i=1}^r n_i (\bar{y}_{i\bullet} - \bar{y}_{..})^2 \quad (2)$$

(SS total) = (SS residual) + (SS treatment)

(SS total) = (SS within) + (SS between)

The corresponding degrees of freedom:

$$n - 1 = (n - r) + (r - 1)$$

with

- $SS \text{ total}$  = total sum of squares

- $SS_{\text{between}}$  = between group of squares, caused by the difference between the groups (treatment effect)
- $SS_{\text{within}}$  = within group of squares, caused by the variation within each group (residual part of the total SS).
- $n = \sum_i n_i$  = total number of observations
- $r$  = number of groups

**Mean sum of squares** = sum of squares divided by the degrees of freedom:

$$\begin{aligned} MS_{\text{total}} &= \frac{SS_{\text{total}}}{n-1} \\ MS_{\text{between}} &= \frac{SS_{\text{between}}}{r-1} \\ MS_{\text{within}} &= \frac{SS_{\text{within}}}{n-r} \end{aligned}$$

**Test statistic** used to test  $H_0 : \mu_1 = \mu_2 = \dots = \mu_r$  versus  $H_1$ : **Not all population averages are the same** is the F-statistic:

$$F = \frac{MS_{\text{between}}}{MS_{\text{within}}} \approx F_{r-1, n-r}$$

**Conclusion:** reject the null hypothesis if F is too big (small p-value)

## 2.2 Example

### Example *Pollution*

What we have just discussed will now be applied to the amount of rain (which is a variable from the data set pollution).

```
glm1 <- lm(Rain ~ regio, data = pollution)
summary(glm1)

##
## Call:
## lm(formula = Rain ~ regio, data = pollution)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -25.5556  -4.6000   0.0417   2.7431  23.0417
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   40.556      3.024  13.411 < 2e-16 ***
## regioN        -5.756      3.825  -1.505 0.138226
## regioNO        1.403      3.546   0.396 0.693954
## regioW       -19.889      4.781  -4.160 0.000115 ***
## regioZO       10.844      5.060   2.143 0.036623 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 9.072 on 54 degrees of freedom
## Multiple R-squared:  0.4239, Adjusted R-squared:  0.3813
## F-statistic: 9.935 on 4 and 54 DF,  p-value: 4.245e-06
```

The average amount of region of the 5 different regions are compared with a **F-test in one-way ANOVA**.  
 $p$ -value (one-sided) = 0.000004 < 0.05

*Conclusion:* reject  $H_0$ . The average amount of rain in these five regions are not the same. The factor **regio** has an effect on the amount of rain. However, we do not know where the differences are.

### 3 One-way ANOVA as a linear model

One-way ANOVA can be seen as a linear model with a continuous response variable and a categorical explanatory variable (with multiple levels).

One-way ANOVA test  $H_0 : \mu_1 = \mu_2 = \dots = \mu_r$  versus  $H_1$ : **Not all the means are the same.**

#### 3.1 Use of treatment coding

We will demonstrate the treatment coding by using the an example.

##### Example *Pollution*

We know, from the descriptive statistics performed earlier, the average amount of rain for the different regions:

Regions	C	N	NO	W	ZO
Average	40.56	34.80	41.96	20.67	51.40

These numbers can be checked by the using the code

Interpretation of the parameter estimates of the output of `summary(glm1)`.

##### Treatment coding

$$Rain = \alpha + \beta_1 D_1 + \beta_2 D_2 + \beta_3 D_3 + \beta_4 D_4 + \epsilon$$

with:

$$D_1 = 1 \text{ if } region = N$$

$$= 0 \text{ else}$$

$$D_2 = 1 \text{ if } region = NO$$

$$= 0 \text{ else}$$

$$D_3 = 1 \text{ if } region = W$$

$$= 0 \text{ else}$$

$$D_4 = 1 \text{ if } region = ZO$$

$$= 0 \text{ else}$$

(3)

The expected values for the different regions:

$$E(Rain_N) = \alpha + \beta_1$$

$$E(Rain_{NO}) = \alpha + \beta_2$$

$$E(Rain_W) = \alpha + \beta_3$$

$$E(Rain_{ZO}) = \alpha + \beta_4$$

$$E(Rain_C) = \alpha$$

(4)

- $\alpha$  = the average amount of rain in  $C$
- $\beta_i$  = difference in average amount of rain in region  $i$  compared to  $C$

→ Hence, region  $C$  can be considered as the reference category.

The output of `summary(glm1)` estimates these parameters:

$$\hat{\alpha} = 40.56$$

$$\hat{\beta}_1 = -5.76$$

$$\hat{\beta}_2 = 1.40$$

$$\hat{\beta}_3 = -19.89$$

$$\hat{\beta}_4 = 10.84$$

Looking back to the means given in the descriptive statistics, we see that the estimated value corresponds to the observed means. For instance for the region  $N$ , we have  $\hat{\alpha} + \hat{\beta}_1 = 40.56 - 5.76 = 34.80$ .

In this model, the null hypothesis  $H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5$  is equivalent to  $H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$ .

In R we can ask for the design matrix (only part of the output is given here)

```
model.matrix(glm1)

##      (Intercept) region regionNO regionW regionZO
## 1             1         1         0         0         0
## 2             1         0         1         0         0
## 3             1         0         1         0         0
## 4             1         0         0         0         1
```

### 3.2 Use of sum coding

When running a linear regression, R will use the *treatment coding* by default. The type of coding can be changed to **sum coding** using (in the function `lm`) the argument `contrasts = list(Name_variable = "contr.sum")` and replacing `Name_variable` by the name of the variable being coded.

#### Example *Pollution*

```
glm2 <- lm(Rain ~ regio, data = pollution, contrasts = list(regio = "contr.sum"))
summary(glm2)

##
## Call:
## lm(formula = Rain ~ regio, data = pollution, contrasts = list(regio = "contr.sum"))
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -25.5556  -4.6000   0.0417   2.7431  23.0417
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   37.876      1.389   27.268 < 2e-16 ***
## regio1         2.679      2.723    0.984  0.3295
## regio2        -3.076      2.285   -1.346  0.1839
## regio3         4.082      1.997    2.044  0.0458 *
## regio4        -17.209      3.187   -5.399 1.53e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 9.072 on 54 degrees of freedom
## Multiple R-squared:  0.4239, Adjusted R-squared:  0.3813
## F-statistic: 9.935 on 4 and 54 DF,  p-value: 4.245e-06
```

Regions	C	N	NO	W	ZO
Average	40.56	34.80	41.96	20.67	51.40

Interpretation of the parameter estimates of the output of `summary(glm2)`:

#### Sum coding



$$Rain = \alpha + \beta_1 D_1 + \beta_2 D_2 + \beta_3 D_3 + \beta_4 D_4 + \epsilon$$

with:

$$\begin{aligned} D_1 &= 1 \text{ if } region = C \\ &= -1 \text{ if } region = ZO \\ &= 0 \text{ else} \\ D_2 &= 1 \text{ if } region = N \\ &= -1 \text{ if } region = ZO \\ &= 0 \text{ else} \\ D_3 &= 1 \text{ if } region = NO \\ &= -1 \text{ if } region = ZO \\ &= 0 \text{ else} \\ D_4 &= 1 \text{ if } region = W \\ &= -1 \text{ if } region = ZO \\ &= 0 \text{ else} \end{aligned} \tag{5}$$

The expectation values for the different regions:

$$\begin{aligned} E(Rain_C) &= \alpha + \beta_1 \\ E(Rain_N) &= \alpha + \beta_2 \\ E(Rain_{NO}) &= \alpha + \beta_3 \\ E(Rain_W) &= \alpha + \beta_4 \\ E(Rain_{ZO}) &= \alpha - \beta_1 - \beta_2 - \beta_3 - \beta_4 \end{aligned} \tag{6}$$

- $\alpha$  = the average of the group averages
- $\beta_i$  = difference in average amount of rain in region  $i$  compared to global average

The output of `summary(glm1)` estimates these parameters:

$$\begin{aligned} \hat{\alpha} &= 37.88 \\ \hat{\beta}_1 &= 2.68 \\ \hat{\beta}_2 &= -3.08 \\ \hat{\beta}_3 &= 4.08 \\ \hat{\beta}_4 &= -17.21 \end{aligned}$$

The expected amount of rain in the region  $ZO$  is:

$$E(Rain_{ZO}) = 37.88 + 2.68 \cdot (-1) - 3.08 \cdot (-1) + 4.08 \cdot (-1) - 17.21 \cdot (-1) = 51.4$$

The model matrix can be returned by using `model.matrix(glm2)`.

```
head(model.matrix(glm2))
```

```
##      (Intercept) regio1 regio2 regio3 regio4
## 1             1      0      1      0      0
## 2             1      0      0      1      0
## 3             1      0      0      1      0
## 4             1     -1     -1     -1     -1
## 5             1      0      0      1      0
## 6             1     -1     -1     -1     -1
```

## 4 Model diagnostics

### 4.1 Assumptions made

ANOVA make the following **assumptions**:

- **Assumption of normality**: Each group sample is drawn from a normally distributed population
- **Assumption of homogeneity of variance**: Different samples have the same variance irrespective if they come from the same population or not
- **Assumption of independence**: The observations between groups should be independent and the observations within each group must be independent.

The first two assumptions in symbolic notation:

$$Y_{1j} \sim N(\mu_1, \sigma^2)$$

$$Y_{2j} \sim N(\mu_2, \sigma^2)$$

...

$$Y_{rj} \sim N(\mu_r, \sigma^2)$$

Besides, it is assumed that the observations are sampled randomly. The presence of outliers can also cause problems and should therefore be checked.

### 4.2 Checking assumptions

**Check the model conditions:**

1. Independent groups?  
→ According to the design of the experiment.
2. Constant within-group variance?  
→ Test  $H_0 : \sigma_1^2 = \sigma_2^2 = \dots = \sigma_r^2$ 
  - Visual check of the boxplot
  - Test for identical variances in a number of groups (e.g., Levene test and Bartlett test)
3. Normal distribution in the groups?  
→ Shapiro-Wilk test per group  
Or  
→ Shapiro-Wilk test of the within-group residuals + histogram of within-group residuals
4. Presence of influential observations  
→ Cook's distance

### 4.3 Check assumption of homogeneity of variance

$$H_0 : \sigma_1^2 = \sigma_2^2 = \dots = \sigma_r^2$$

versus

$$H_1 : \text{not all } \sigma_i^2 \text{ are equal } (i = 1, 2, \dots, r)$$

We can use the Levene's test to test the homogeneity of variance assumption.

In R, this test can be performed by the function `leveneTest` from the package `car`.

```
library(car)
leveneTest(rain ~ region)
```

```
## Levene's Test for Homogeneity of Variance (center = median)
##      Df F value Pr(>F)
## group  4  1.6305 0.1799
##      54
```

We can assume homogeneity of variances since  $p\text{-value} > 0.05$ .

### Remark 1:

If the Levene's test is rejected, be aware that there exists some robustness. If the variances are not too unequal, we can still make use of the ANOVA F-test.

*Rule of thumb:*  $\frac{\max(\sigma_1, \sigma_2, \dots, \sigma_r)}{\min(\sigma_1, \sigma_2, \dots, \sigma_r)} \leq 5$

### Remark 2:

In case there is no homogeneity of variances, a modification of the F-test can be used:

```
oneway.test(Rain ~ regio, data = pollution, var.equal = FALSE)
```

```
##
## One-way analysis of means (not assuming equal variances)
##
## data: Rain and regio
## F = 12.07, num df = 4.000, denom df = 14.202, p-value = 0.0001751
```

## 4.4 Check assumption of normality

We assume

$$Y_{1j} \sim N(\mu_1, \sigma^2) \Rightarrow Y_{1j} - \mu_1 \sim N(0, \sigma^2)$$

$$Y_{2j} \sim N(\mu_2, \sigma^2) \Rightarrow Y_{2j} - \mu_2 \sim N(0, \sigma^2)$$

...

$$Y_{rj} \sim N(\mu_r, \sigma^2) \Rightarrow Y_{rj} - \mu_r \sim N(0, \sigma^2)$$

Hence, we can check normality for the within-group residuals!

1. What are the residuals in our example?

```
by_region <- group_by(pollution, regio)
summarise(by_region, n = n(), mean = mean(Rain), sd = sd(Rain), median = median(Rain),
          min = min(Rain), max = max(Rain))
```

```
## # A tibble: 5 x 7
##   regio      n mean    sd median   min   max
##   <fct> <int> <dbl> <dbl> <dbl> <int> <int>
## 1 C      9  40.6  11.9    45    15    54
## 2 N     15  34.8   3.10    35    30    40
## 3 NO    24  42.0   9.99    42    25    65
## 4 W      6  20.7  12.2   15.5    10    37
## 5 ZO      5  51.4   5.86    52    45    60
```

Look at the first 5 observations of the pollution data set

```
head(pollution, n = 5)
```

```
##      city regio JanT JulT Hum Rain Mortality Educ density  NW
## 1   Akron     N   27   71  59   36   921.87  11.4   3243  8.8
## 2  Albany    NO   23   72  57   35   997.87  11.0   4281  3.5
## 3 Allentown  NO   29   74  54   44   962.35   9.8   4260  0.8
## 4  Atlanta   ZO   45   79  56   47   982.29  11.1   3125 27.1
## 5 Baltimore NO   35   77  55   43  1071.29   9.6   6441 24.4
```

Compute the corresponding residuals by yourself

City	Residual
Akron	
Albany	
Allentown	

City	Residual
Atlanta	
Baltimore	

Compare your results with the residuals obtained in R

```
glm1 <- lm(Rain ~ regio, data = pollution)
combine <- data.frame(city = pollution$city, region = pollution$regio,
                      residuals = glm1$residuals)
head(combine, n = 5)
```

```
##      city region residuals
## 1   Akron      N  1.200000
## 2  Albany     NO -6.958333
## 3 Allentown    NO  2.041667
## 4  Atlanta     ZO -4.400000
## 5 Baltimore    NO  1.041667
```

## 2. How to test normality of these residuals in R?

Test normality of the within-group residuals:

```
shapiro.test(glm1$residuals)
```

```
##
##  Shapiro-Wilk normality test
##
## data:  glm1$residuals
## W = 0.94388, p-value = 0.008825
```

```
hist(glm1$residuals)
```



**Remark:**

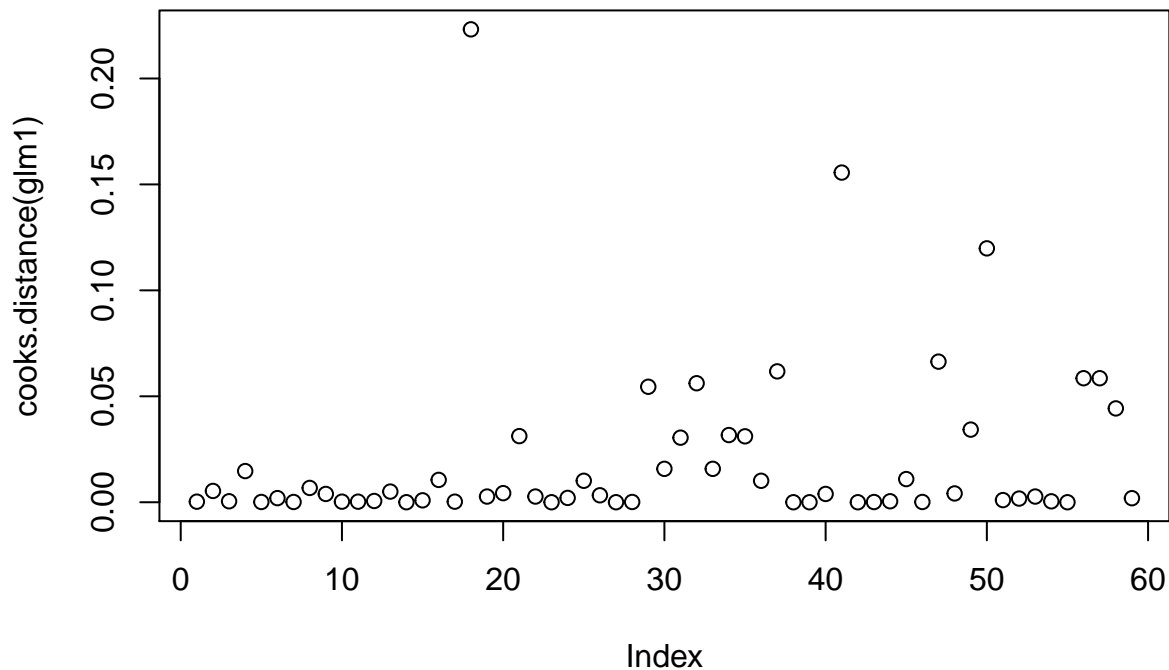
- When normality is rejected by the Shapiro-Wilk test, we still can interpret the ANOVA F statistic as long as the residuals are symmetric distributed (seen from the histogram).
- In case of asymmetric distribution of residuals, we can transform the response variable  $Y$ .
- In case of asymmetry, we can always use the non-parametric version: the Kruskal-Wallis test.

## 4.5 Checking for influential observations

As in regression analysis, we use the **Cook's distance** to check for influential observations.

We plot the Cook's distance versus observation number

```
plot(cooks.distance(glm1))
```



There are no influential observations.

## 5 Pairwise comparisons of treatment effects

### 5.1 Questions post-hoc

We rejected the null hypothesis  $H_0 : \mu_N = \mu_{NO} = \mu_W = \mu_{ZO} = \mu_C$  with  $\mu_i$  the average amount of rain in region  $i$ . The rejection was made based on the results of a F-test in one-way ANOVA. (see earlier).

We now raise some questions post-hoc:

1. Is there a pair of regions with the same average amount of rain?
2. Test all pairs of regions on the same average amount of rain.

If you want to do multiple comparisons in R, you have to first use the function `aov` and store the results in an object.

```
poll.aov1 <- aov(Rain ~ regio, data = pollution)
summary(poll.aov1)
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## regio      4   3270   817.6    9.935 4.24e-06 ***
## Residuals 54   4444    82.3
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

## 5.2 Planned comparison of means in ANOVA

When performing an ANOVA test and a significant result is obtained, then the null hypothesis of equal means is rejected. However, an ANOVA test cannot tell which group differs. To address this problem, the **Least Significant Difference** (LSD-method) method can be used to test a planned comparison (between two groups). In the LSD test, the mean of one group is compared to the mean of another group. The LSD test is basically a t-test for two means in ANOVA. The main difference with a regular t-test is that the standard deviation estimate is based on the observations of all the groups. In a regular t-test, only the observations from the two groups under consideration are used to estimate the standard deviation.

Consider  $r$  groups with population means  $\mu_1, \mu_2, \dots, \mu_r$ . Assume an ANOVA test is performed and suggest rejection of the null hypothesis, that is, the difference between the group means is significant. We now want to check whether the mean of group 1 and group 2 are significantly different from each other. This is checked with a LSD test.

### Test problem:

$H_0 : \mu_1 = \mu_2$  (group 1 and group 2 have the same mean)

$H_1 : \mu_1 \neq \mu_2$  (group 1 and group 2 have a different mean)

### Test statistic:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{s^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \sim t_{n-r} \text{ (distribution under } H_0 \text{)}$$

with  $s^2 = MS \text{ within} = \frac{SS \text{ within}}{n-r} = \frac{\sum_{i=1}^r \sum_{j=1}^{n_j} (Y_{ij} - \bar{Y}_{i\bullet})^2}{n-r}$  the pooled estimator of the variance.

### Conclusion:

Reject the null hypothesis and accept a difference in effect between the two groups if the test statistic  $t$  is too big or too small (two-sided  $p$ -value).

This test method is also known as the LSD-method (Least Significant Difference method).

To compare two groups in ANOVA:

- In case there is homogeneity of variances, use the LSD method in ANOVA (with  $n - r$  degrees of freedom, see above).
- In case there is no homogeneity of variances, use the two-sample t-test for independent groups (with  $n_1 + n_2 - 2$  degrees of freedom).

The usual t-test for two groups estimates the within-group variance  $\sigma^2$  based on the pooled sample variance of the two groups. For each pair of groups, one uses a different variance estimator, which is not logical because all groups have the same variance. The t-test above for two groups in ANOVA (thus the LSD test) uses information from all groups – the pooled sample variance based on the  $r$  groups – to estimate  $\sigma^2$ . Thus this test will make less wrong decisions. On the other hand, if there are indications that not all groups have the same variance, one should prefer the usual t-test for two groups.

### Example *Pollution in R*

We want to test whether the mean amount of rain in region  $C$  and region  $NO$  are the same or not.

### Statement hypothesis:

$H_0 : \mu_{NO} = \mu_C$  versus  $H_1 : \mu_{NO} \neq \mu_C$

```
pairwise.t.test(rain, region, p.adjust.method = "none")
```

```
##
## Pairwise comparisons using t tests with pooled SD
##
## data:  rain and region
##
##      C      N      NO      W
## N 0.13823 -      -      -
```

```
## NO 0.69395 0.02000 - -
## W 0.00011 0.00214 3.9e-06 -
## ZO 0.03662 0.00082 0.03887 7.5e-07
##
## P value adjustment method: none
```

### Conclusion:

There is no significant difference in the average amount of rain between the region *NO* and *C*.

## 5.3 Multiple comparisons

### 5.3.1 The problem of multiple comparisons

#### A planned test versus multiple test

If the groups, one wants to compare, are determined before one looks at the data, then the LSD test is suitable.

#### Risk on test differences which in reality are not present

If one compares each pair of means using the LSD test with significance level  $\alpha = 0.05$ , there is a *high probability that one falsely reports differences*.

(remember:  $\alpha = P(\text{reject } H_0 | H_0 \text{ true})$ )

- Even with a small number of groups, many pairs are possible, i.e.

$$\binom{r}{2} = \frac{r(r-1)}{2} \quad (7)$$

- A t-test for two means ( $\mu_1 = \mu_2$  versus  $\mu_1 \neq \mu_2$ ) with significance level  $\alpha = 0.05$  will report a false difference in 5% of applications (or 1 on 20 samples) on populations with the same mean. (e.g. a clinical study placebo/medication: a medicine without effect will be concluded as effective in 5% of tests, due to the coincidental structure of the sample.)

**Example:** ANOVA with  $r = 10$  groups, gives  $\frac{10 \cdot 9}{2} = 45$  pairs. If there are no population differences, and one tests each pair with a LSD test on significance level  $\alpha = 0.05$ , one can expect that in these samples,  $5\% \cdot 45 = 2$  pairs will be considered different, while in the populations they are not. (Roughly, because the successive tests are not statistical independent experiments.)

**Problem:** find a procedure that, if there are no differences, using multiple comparisons for all pairs, keeps the overall probability of falsely reporting at least one difference, less than a certain chosen significance level  $\alpha$ .

### 5.3.2 An overview of multiple comparisons procedures

1. **Tukey HSD (honestly significant differences) test for multiple pairs: all pairwise comparisons**
  - Test similarity of pairs:  $H_0 : \mu_i = \mu_j$  (all pairs are equal)  
Reports the differences pairwise  $\mu_i \neq \mu_j$
  - Advantage: The overall significance level is exactly  $\alpha$
2. **Bonferroni method for multiple tests**
  - Test similarity of pairs:  $H_0 : \mu_i = \mu_j$
  - Per pair t-test with reduced significance level:  $\alpha^* = \frac{\alpha}{\text{number of pairs}}$
  - Disadvantage: overall significance level  $\leq \alpha$ , even  $< \alpha$  (test is conservative for  $H_0$ )
  - Advantage: you can use it with a small number of groups
3. **Scheffé multiple contrasts test: all linear contrasts**
  - $H_0 : c_1\mu_1 + c_2\mu_2 + \dots + c_r\mu_r = 0$  (given that  $\sum_{i=1}^r c_i = 0$ ) (all linear contrasts 0)
  - Disadvantage: overall significance level  $\leq \alpha$  (test sometimes conservative for  $H_0$ )
4. **Holm's step-wise correction for multiple tests**



- Test similarity of pairs:  $H_0 : \mu_i = \mu_j$
- The Holm adjustment sequentially compares the lowest  $p$ -value with a type I error rate that is reduced for each consecutive test. This method is generally considered superior to the Bonferroni adjustment.

## In R

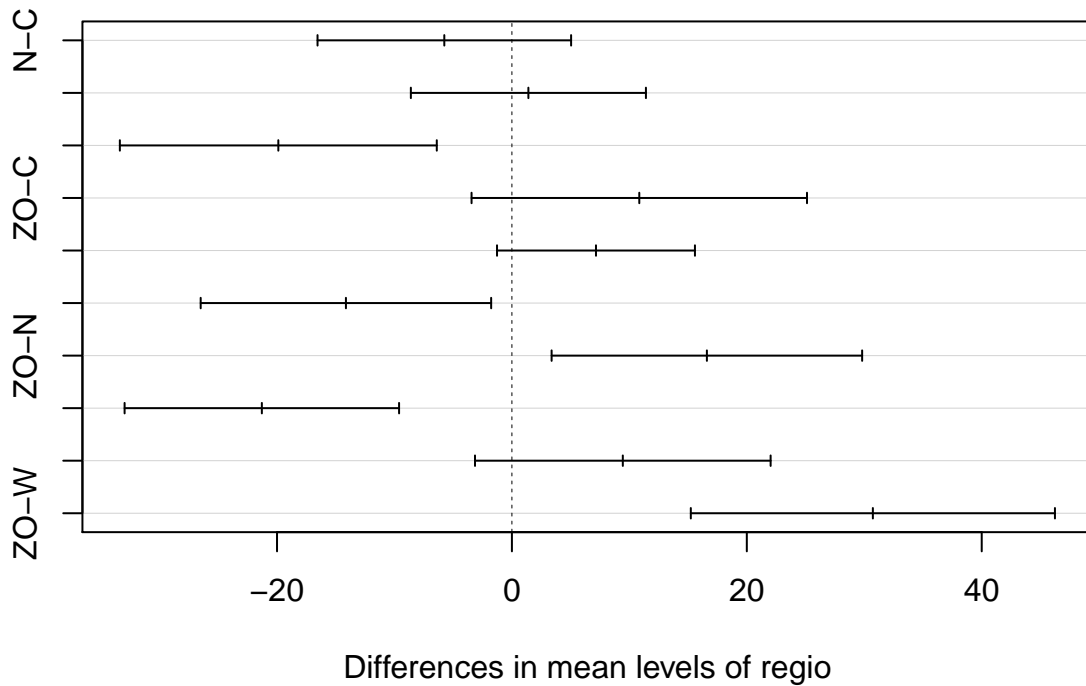
Tukey HSD test

```
diffs <- TukeyHSD(poll.aov1, which = "region", conf.level = 0.95)
diffs
```

```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = Rain ~ regio, data = pollution)
##
## $regio
##          diff          lwr          upr      p adj
## N-C      -5.755556 -16.550086  5.038975 0.5639861
## NO-C       1.402778  -8.604020 11.409575 0.9946762
## W-C     -19.888889 -33.382052 -6.395726 0.0010501
## ZO-C      10.844444  -3.435377 25.124266 0.2172902
## NO-N       7.158333  -1.268144 15.584811 0.1316844
## W-N     -14.133333 -26.500021 -1.766645 0.0174879
## ZO-N      16.600000  3.379454 29.820546 0.0070724
## W-NO     -21.291667 -32.977089 -9.606245 0.0000371
## ZO-NO      9.441667  -3.143918 22.027251 0.2278970
## ZO-W      30.733333 15.230869 46.235797 0.0000073
```

```
plot(diffs)
```

## 95% family-wise confidence level



If the value 0 is not included in the 95 % CI then there is a significant difference between the average values of the two groups.

The adjusted  $p$  - value returned by the `aov` function depends on assumptions of the residuals. For the  $p$  - value to be correct, these residuals need to be independent, normally distributed, and have constant variance. In the following section, we see a non-parametric function that does not require the normality assumption.

### Remark 1:

You can use the `scheffe.test` function (from package `agricolae`) for Scheffé multiple comparisons and the function `pairwise.t.test` for the Bonferroni multiple comparisons. There exist also a HSD function which is similar to Tukey.

```
library(agricolae)
```

```
poll.aov1 <- aov(Rain ~ regio, data = pollution)
```

another HSD test

```
HSD.test(poll.aov1, "regio", group = FALSE)$comparison
```

##	difference	pvalue	signif.	LCL	UCL
## C - N	5.755556	0.5640		-5.038975	16.550086
## C - NO	-1.402778	0.9947		-11.409575	8.604020
## C - W	19.888889	0.0011	**	6.395726	33.382052
## C - ZO	-10.844444	0.2173		-25.124266	3.435377
## N - NO	-7.158333	0.1317		-15.584811	1.268144
## N - W	14.133333	0.0175	*	1.766645	26.500021
## N - ZO	-16.600000	0.0071	**	-29.820546	-3.379454

```
## NO - W    21.291667 0.0000    ***    9.606245  32.977089
## NO - ZO   -9.441667 0.2279         -22.027251   3.143918
## W - ZO   -30.733333 0.0000    ***   -46.235797 -15.230869
```

Scheffé multiple contrasts test

```
scheffe.test(poll.aov1, "region", group = FALSE)$comparison
```

```
##          Difference pvalue sig          LCL          UCL
## C - N      5.755556 0.6883         -6.4436228  17.954734
## C - NO     -1.402778 0.9970        -12.7117188   9.906163
## C - W     19.888889 0.0042    **    4.6399160  35.137862
## C - ZO    -10.844444 0.3439        -26.9824405   5.293552
## N - NO     -7.158333 0.2344        -16.6813139   2.364647
## N - W     14.133333 0.0461    *    0.1574188  28.109248
## N - ZO    -16.600000 0.0215    *   -31.5408811  -1.659119
## NO - W     21.291667 0.0002    ***    8.0856687  34.497665
## NO - ZO     -9.441667 0.3565        -23.6649617   4.781628
## W - ZO    -30.733333 0.0000    ***   -48.2530694 -13.213597
```

Bonferroni method

```
pairwise.t.test(rain, region, p.adj="bonferroni")
```

```
##
## Pairwise comparisons using t tests with pooled SD
##
## data:  rain and region
##
##      C      N      NO      W
## N  1.0000 -      -      -
## NO 1.0000 0.2000 -      -
## W  0.0011 0.0214 3.9e-05 -
## ZO 0.3662 0.0082 0.3887 7.5e-06
##
## P value adjustment method: bonferroni
```

### Remark 2: Homogeneous groups

Scheffé, homogeneous groups

```
Hgroups.scheffe <- scheffe.test(poll.aov1, "region", group = TRUE)
Hgroups.scheffe$groups
```

```
##          Rain groups
## ZO 51.40000      a
## NO 41.95833     ab
## C  40.55556     ab
## N  34.80000      b
## W  20.66667      c
```

### Remark 3: Holm's approach

*Illustration of the Holm's approach*

1. We compute the non-adjusted  $p$  – values for every hypothesis (In this example: 10 different pairs so 10  $p$  – values)

```
pairwise.t.test(rain, region, p.adj = "none")
```

```
##
```

```
## Pairwise comparisons using t tests with pooled SD
```

```
##
```

```
## data: rain and region
```

```
##
```

```
##      C      N      NO      W
## N  0.13823 -      -      -
## NO 0.69395 0.02000 -      -
## W   0.00011 0.00214 3.9e-06 -
## ZO 0.03662 0.00082 0.03887 7.5e-07
```

```
##
```

```
## P value adjustment method: none
```

2. Compare the smallest  $p$ -value with  $\frac{0.05}{10} = 0.005$   
This is  $7 \cdot 10^{-7} < 0.005$ . Hence this null hypothesis is rejected. There is a significant difference in average rain between  $ZO$  and  $W$ .
3. Compare the second smallest  $p$ -value to  $\frac{0.05}{9} = 0.0056$ .  
This is  $3.9 \cdot 10^{-6} < 0.0056$ . Hence this null hypothesis is rejected. There is a significant difference in average rain between  $W$  and  $NO$ .
4. Compare the third smallest  $p$ -value to  $\frac{0.05}{8} = 0.0063$ .  
This is  $0.00011 < 0.0063$ . Hence this null hypothesis is rejected. There is a significant difference in average rain between  $W$  and  $N$ .
5. Compare the fourth smallest  $p$ -value to  $\frac{0.05}{7} = 0.0071$   
This is  $0.0008 < 0.0071$ . Hence the corresponding null hypothesis is rejected.
6. Compare the fifth smallest  $p$ -value to  $\frac{0.05}{6} = 0.0083$   
This is  $0.002 < 0.0083$ . Hence the corresponding null hypothesis is rejected.
7. Compare the sixth smallest  $p$ -value to  $\frac{0.05}{5} = 0.01$   
This is  $0.02 > 0.01$ . Hence the corresponding null hypothesis is not rejected. As soon as that happens, you stop, and therefore, also fail to reject the remaining hypothesis

### Holm's approach in R

```
pairwise.t.test(rain, region, p.adj = "holm")
```

```
##
```

```
## Pairwise comparisons using t tests with pooled SD
```

```
##
```

```
## data: rain and region
```

```
##
```

```
##      C      N      NO      W
## N  0.27645 -      -      -
## NO 0.69395 0.10001 -      -
## W   0.00092 0.01283 3.5e-05 -
## ZO 0.14649 0.00577 0.14649 7.5e-06
```

```
##
```

```
## P value adjustment method: holm
```

**Remark 4:** General remark about the use of the `aov` function in R

Only factors can be used in ANOVA. The `aov` function really needs the explanatory variables to be a factor. An error is returned when an explanatory variable `var1` is not a factor. This can be solved by making variable `var1` as a factor `var1.f` (use function `as.factor()`).

```
var1.f <- as.factor(var1)
```

## 6 Extensions

### 6.1 The Kruskal-Wallis test

The **Kruskal-Wallis test** is a **non-parametric version of one-way analysis of variance**. The assumption underlying this test is that the measurements come from a continuous distribution, but not necessarily a normal distribution. The test is based on an analysis of variance using the ranks of the data values, not the data values.

**Statement of hypothesis:**

$H_0$  : The location parameters of the distribution of  $X$  are the same in each group.

versus

$H_1$  : The location parameters differ in at least one group.

**Test statistic:**

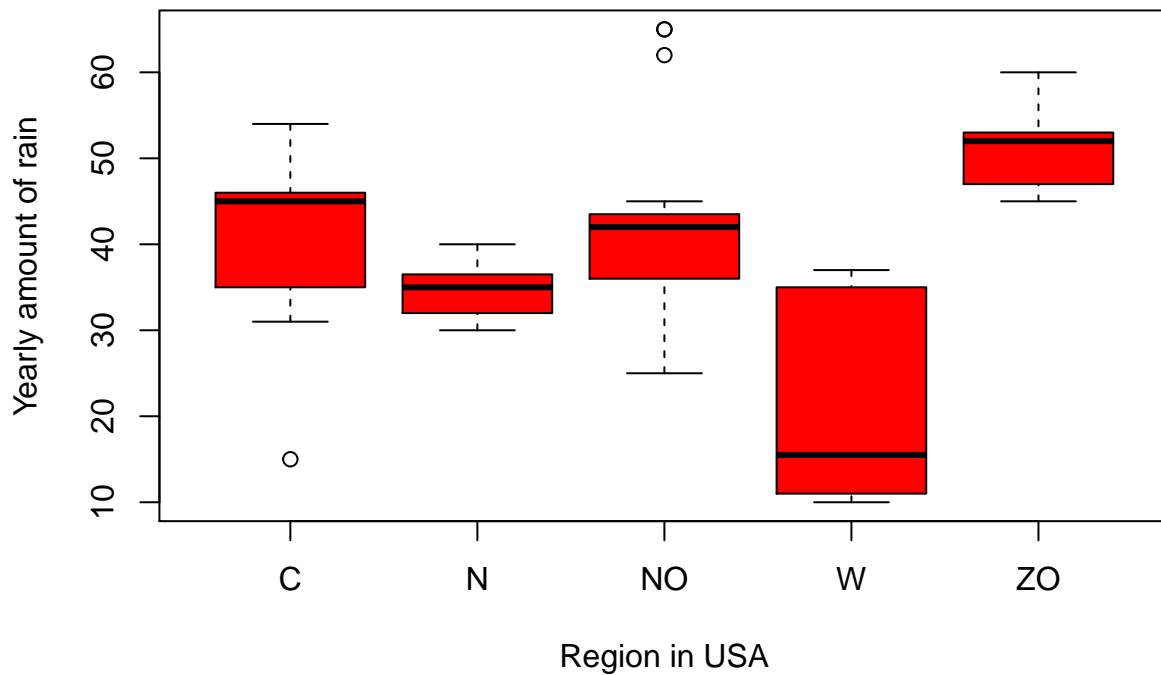
1. We sort our data from small to large values for the variable *Rain*.

```
subpol <- pollution[,c("Rain", "region")]
library("doBy")
sortsubpol <- orderBy( ~ Rain, data = subpol)
head(sortsubpol)
```

```
##      Rain region
## 47    10      W
## 29    11      W
## 49    13      W
## 18    15      C
## 48    18      W
## 34    25     NO
```

2. We associate ranks to the observations. The smallest observation gets rank 1, the largest one gets rank 59. We then compute the total number of ranks within each region.

Region	number of observations	Total number of ranks	Average rank
<i>C</i>	9	329	36.5
<i>N</i>	15	304	20.3
<i>NO</i>	24	818.5	34.1
<i>W</i>	6	57.5	9.6
<i>ZO</i>	5	261	52.2



3. The Kruskal-Wallis statistic is based on the ranks

Compute the Kruskal-Wallis statistic in R:

```
kruskal.test(rain ~ region)
```

```
##
## Kruskal-Wallis rank sum test
##
## data: rain by region
## Kruskal-Wallis chi-squared = 24.397, df = 4, p-value = 6.649e-05
```

### Conclusion:

Since  $p\text{-value} < 0.05$ , we reject  $H_0$ . The location parameters are not the same in all the regions. In at least one region, there is a shift in location parameter.

4. Multiple comparisons

When you install the R package `pgirmess`, you have the possibility to ask for multiple comparisons:

```
library(pgirmess)
```

```
kruskalmc(rain, region)
```

```
## Multiple comparison test after Kruskal-Wallis
## p.value: 0.05
## Comparisons
##      obs.dif critical.dif difference
## C-N    16.28889    20.32813     FALSE
## C-NO    2.451389    18.84468     FALSE
## C-W    26.97222    25.41016      TRUE
## C-ZO    15.64444    26.89159     FALSE
```

```
## N-NO 13.837500    15.86864    FALSE
## N-W  10.683333    23.28880    FALSE
## N-ZO 31.933333    24.89677     TRUE
## NO-W 24.520833    22.00584     TRUE
## NO-ZO 18.095833    23.70102    FALSE
## W-ZO 42.616667    29.19405     TRUE
```

## 7 Two-way ANOVA

### 7.1 Introduction

Consider an experiment involving two fixed-effect factors.

#### Example *Diet*

24 men, each weighing about 20 kg too much, are accidentally spread over 12 treatments (3 levels of jogging, 4 levels of diet). There are two men for each treatment. Each man used the same number of calories per day, the diet differs only in the amounts of protein, fat and carbohydrates. At the end of the experiment, the men are weighted again and their losses are calculated. At the end of the experiment, the men are weighted again and their losses are calculated.

<i>Weight loss</i>		<b>Diet</b>			
		Normal	Protein	Fat	Carbohydrates
<b>Jogging</b>	0 km	8.5	15.5	8.5	15.5
		11.5	16.5	7.5	13.5
	1 km	14	20	13	21
		16	23	11	18
	2 km	24.5	27	22	24.5
		19.5	24	27	27.5

Import the data set *diet.txt* as `diet_df` in R.

```
head(diet_df)
```

```
##   LOSS JOGGING   DIET
## 1  8.5      0km normal
## 2 11.5      0km normal
## 3 15.5      0km protein
## 4 16.5      0km protein
## 5  8.5      0km   fat
## 6  7.5      0km   fat
```

#### *Descriptive statistics*

We can ask for some description statistics for the treatment:

```
names(diet_df)
```

```
## [1] "LOSS"      "JOGGING"   "DIET"
```

```
by_diet_jogging <- group_by(diet_df, DIET, JOGGING)
summarise(by_diet_jogging, Avgloss = mean(LOSS), SDloss=sd(LOSS), number=n())
```

```
## # A tibble: 12 x 5
## # Groups:   DIET [4]
##   DIET      JOGGING Avgloss SDloss number
##   <fct>    <fct>      <dbl> <dbl> <int>
## 1 carbo    0km         14.5  1.41     2
## 2 carbo    1km         19.5  2.12     2
## 3 carbo    2km         26     2.12     2
## 4 fat      0km          8     0.707    2
## 5 fat      1km         12     1.41     2
## 6 fat      2km        24.5  3.54     2
## 7 normal   0km         10     2.12     2
## 8 normal   1km         15     1.41     2
## 9 normal   2km         22     3.54     2
## 10 protein 0km         16     0.707    2
## 11 protein 1km        21.5  2.12     2
## 12 protein 2km        25.5  2.12     2
```

If we want to see the descriptive statistics for jogging and diet separately:

- descriptive statistics by diet

```
by_diet <- group_by(diet_df, DIET)
summarise(by_diet, Avgloss = mean(LOSS), SDloss=sd(LOSS), number=n())
```

```
## # A tibble: 4 x 4
##   DIET      Avgloss SDloss number
##   <fct>      <dbl> <dbl> <int>
## 1 carbo      20     5.37     6
## 2 fat       14.8    7.89     6
## 3 normal    15.7    5.73     6
## 4 protein    21     4.48     6
```

- descriptive statistics by jogging

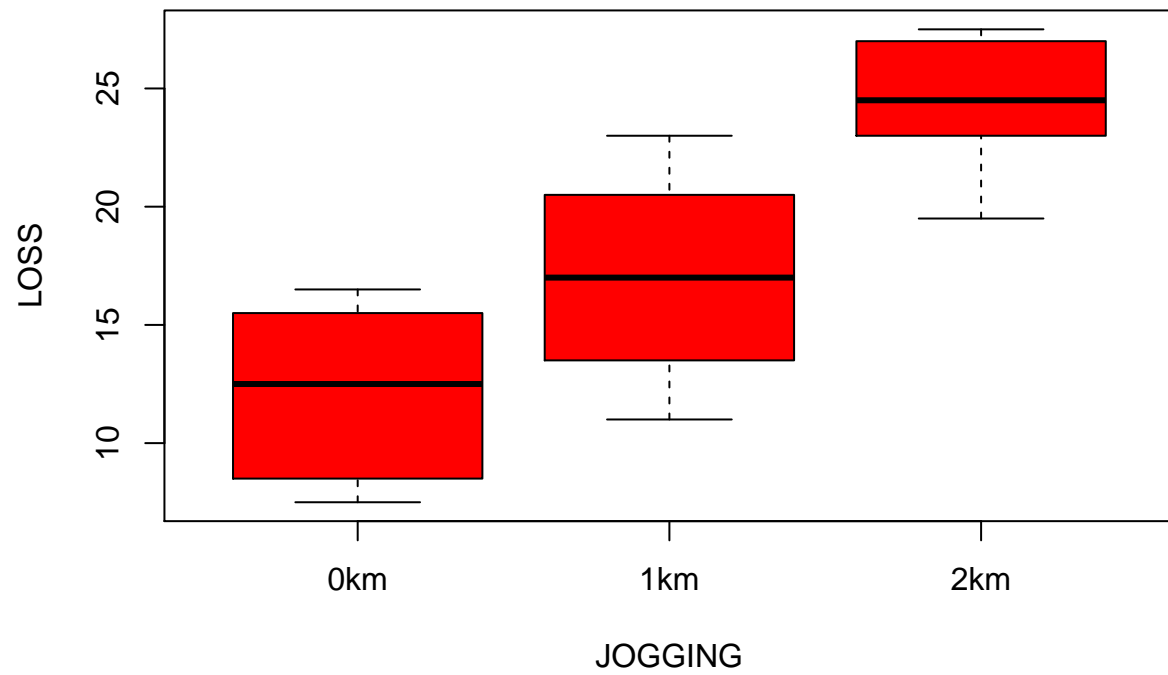
```
by_jogging <- group_by(diet_df, JOGGING)
summarise(by_jogging, Avgloss = mean(LOSS), SDloss=sd(LOSS), number=n())
```

```
## # A tibble: 3 x 4
##   JOGGING Avgloss SDloss number
##   <fct>      <dbl> <dbl> <int>
## 1 0km      12.1    3.62     8
## 2 1km      17     4.21     8
## 3 2km      24.5    2.75     8
```

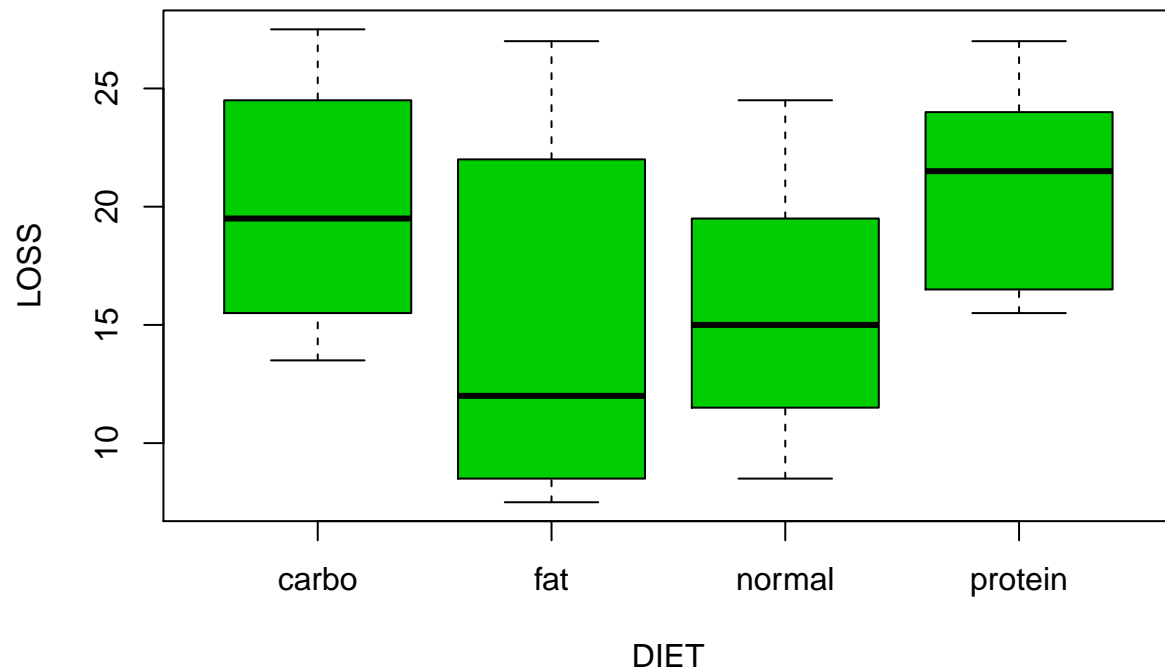
Some univariate box plots

```
boxplot(LOSS ~ JOGGING, col=2, data=diet_df)
```





```
boxplot(LOSS ~ DIET, col=3, data=diet_df)
```

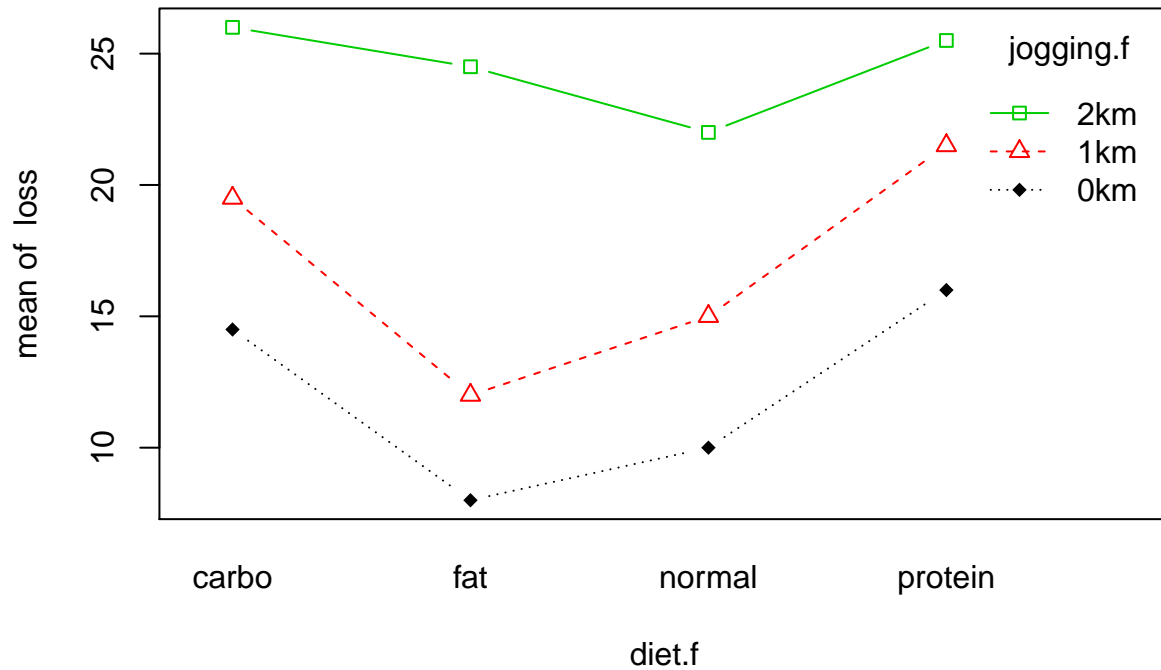


When we only consider one factor at a time, we miss the ‘joint’ effect. Such a joint effect is called **interaction**. We can also ask for one interaction plot with the function `interaction.plot` from the package `stats`. This function plots the mean (or other summary) of the response for two-way combinations of factors, thereby illustrating possible interactions.

```
?interaction.plot
```

Visualization of the mean values for every combination of the factors `diet` and `jogging`

```
diet.f <- as.factor(diet_df$DIET)
jogging.f <- as.factor(diet_df$JOGGING)
loss <- diet_df$LOSS
interaction.plot(diet.f, jogging.f, loss, type = "b", pch = c(18, 24, 22), col = c(1, 2, 3))
```



## 7.2 Two-way ANOVA model

---

Regression model

---

Model without interaction:  $Y = \alpha + \beta_1 x_1 + \beta_2 x_2$

Model with interaction:  $Y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} (x_1 \cdot x_2)$

---



---

ANOVA model

---

Main effects model:  $Y_{ij} = \mu_{\bullet\bullet} + \alpha_i + \beta_j$

Model with interaction:  $Y_{ij} = \mu_{\bullet\bullet} + \alpha_i + \beta_j + (\alpha\beta)_{ij}$

---

In the two-way ANOVA model with interaction:

$$Y_{ijk} = \mu_{\bullet\bullet} + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk}$$

with  $i = 1, 2, \dots, I$ ;  $j = 1, 2, \dots, J$  and  $k = 1, 2, \dots, K$

### Example *Diet*

The weight loss for an individual can be written as the sum of

- The overall mean loss ( $\mu_{\bullet\bullet}$ )
- A part depending on the jogging level ( $\alpha_i$ )
- A part depending on the diet type ( $\beta_i$ )
- A part depending on the interaction between jogging and diet type ( $(\alpha\beta)_{ij}$ )
- An error term  $\varepsilon_{ijk}$  which we assume to be independent and normally distributed ( $\varepsilon_{ijk} \sim N(0, \sigma^2)$ )

## 7.3 Strategy for the analysis of two-way ANOVA studies

**Step 1:** Test whether the interaction is significant

Yes → If yes, go to step 2A

No → If no, go to step 2B

**Step 2A:** The interaction term is significant

- Check the diagnostics.
- Use pairwise comparisons on interaction effect.

**Step 2B:** Only main effects are important. Drop the interaction term and refit the model.

- Check the diagnostics.
- Use pairwise comparisons on the main effects.

### 7.3.1 Example in R

#### Example *Diet* in R

We always use **sum coding** when working with two-way ANOVA.

```
diet.aov1 <- aov(LOSS ~ JOGGING + DIET + JOGGING*DIET, data = diet_df,
               contrasts = list(JOGGING = "contr.sum", DIET = "contr.sum"))
summary(diet.aov1)
```

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## JOGGING        2  621.7   310.87   68.450 2.74e-07 ***
## DIET           3  170.5    56.82   12.511 0.000528 ***
## JOGGING:DIET    6   43.9     7.32    1.612 0.226638
## Residuals     12   54.5     4.54
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

We have an F-test for each effect in our model.

**Step 1:** Check whether the interaction term is significant.

The interaction term is not significant ( $p$ -value = 0.23).

This means that the curves in the `interaction.plot` are parallel.

**Step 2B:** Hence we can drop the interaction term from our model and rerun the model.

Drop interaction term and refit model:

```
diet.aov2 <- aov(LOSS ~ JOGGING + DIET, data = diet_df,
               contrasts = list(JOGGING = "contr.sum", DIET = "contr.sum"))
summary(diet.aov2)
```

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## JOGGING        2  621.7   310.87   56.86 1.66e-08 ***
## DIET           3  170.5    56.82   10.39 0.00034 ***
## Residuals     18   98.4     5.47
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

We now have an additive model with 2 significant main factors: *JOGGING* and *DIET*.

#### Remark:

Since ANOVA can be seen as a linear model, we can also use the `lm` function in R.

Here we give the results for the full model (comparable with the previous results).

```
diet.lm <- lm(LOSS ~ JOGGING + DIET + JOGGING * DIET, data = diet_df,
             contrasts = list(JOGGING = "contr.sum", DIET = "contr.sum"))
# To see the ANOVA table
Anova(diet.lm, type = "III") # function from package 'car'
```

```
## Anova Table (Type III tests)
##
## Response: LOSS
##              Sum Sq Df    F value    Pr(>F)
## (Intercept)  7668.4  1 1688.4495 2.795e-14 ***
## JOGGING       621.7  2   68.4495 2.740e-07 ***
## DIET          170.5  3   12.5107 0.0005277 ***
## JOGGING:DIET   43.9  6    1.6116 0.2266382
## Residuals     54.5 12
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

## 7.4 Diagnostics

### 7.4.1 Checking homogeneity of variances

When the interaction term is significant, then we have to check

$$H_0 : \sigma_{11}^2 = \sigma_{12}^2 = \dots \sigma_{IJ}^2$$

Since the interaction term is not significant (in the *Diet* example), it suffices to check homogeneity of variances for *JOGGING* and *DIET* separately.

- **For *JOGGING*:**  $H_0 : \sigma_{0km}^2 = \sigma_{1km}^2 = \sigma_{2km}^2$   
Test homogeneity of variances for *JOGGING*:

```
leveneTest(LOSS ~ JOGGING, data = diet_df) # This is a function from the package 'car'
```

```
## Levene's Test for Homogeneity of Variance (center = median)
##      Df F value Pr(>F)
## group  2  1.6545 0.2151
##      21
```

- **For *DIET*:**  $H_0 : \sigma_{normal}^2 = \sigma_{protein}^2 = \sigma_{fat}^2 = \sigma_{carbo}^2$   
Test homogeneity of variances for *DIET*:

```
leveneTest(LOSS ~ DIET, data = diet_df)
```

```
## Levene's Test for Homogeneity of Variance (center = median)
##      Df F value Pr(>F)
## group  3  0.3822 0.7669
##      20
```

#### Remark:

1. In case Levene's test is rejected, we can also use the rule of thumb:

$$\frac{\max(\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2)}{\min(\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2)} \leq 5. \text{ In case the rule of thumb is satisfied, an ANOVA F-test can be used.}$$

```
diet_df %>% group_by(JOGGING, DIET) %>%
  summarise(mean = mean(LOSS, na.rm = T), var = var(LOSS, na.rm = T),
            n = n()) %>% select(DIET, JOGGING, mean, var, n)
```

```
## # A tibble: 12 x 5
## # Groups:   JOGGING [3]
##   DIET    JOGGING mean  var    n
```

```
##      <fct>   <fct>   <dbl> <dbl> <int>
##  1 carbo    0km      14.5    2      2
##  2 fat      0km       8      0.5    2
##  3 normal   0km      10      4.5    2
##  4 protein  0km      16      0.5    2
##  5 carbo    1km      19.5    4.5    2
##  6 fat      1km      12      2      2
##  7 normal   1km      15      2      2
##  8 protein  1km      21.5    4.5    2
##  9 carbo    2km      26      4.5    2
## 10 fat      2km      24.5   12.5    2
## 11 normal   2km      22      12.5    2
## 12 protein  2km      25.5    4.5    2
```

2. If there is an indication of unequal variances (based on the above tests), you can use heteroscedastic consistent covariance matrices.

(Here, this is not the case. We just illustrate here how to work. )

```
# Robust analysis
Anova(diet.aov2, type = "III", white.adjust = "hc3")

## Analysis of Deviance Table (Type III tests)
##
## Response: LOSS
##              Df          F    Pr(>F)
## (Intercept)  1 1051.8855 < 2.2e-16 ***
## JOGGING      2   42.2050 1.601e-07 ***
## DIET         3    7.9432 0.001395 **
## Residuals    18
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

This robust analysis also indicates that *JOGGING* and *DIET* are significant.

## 7.4.2 Checking normality of residuals

Test normality of the within-cell residuals

```
shapiro.test(diet.aov2$residuals)

##
## Shapiro-Wilk normality test
##
## data:  diet.aov2$residuals
## W = 0.97129, p-value = 0.699
hist(diet.aov2$residuals)
```

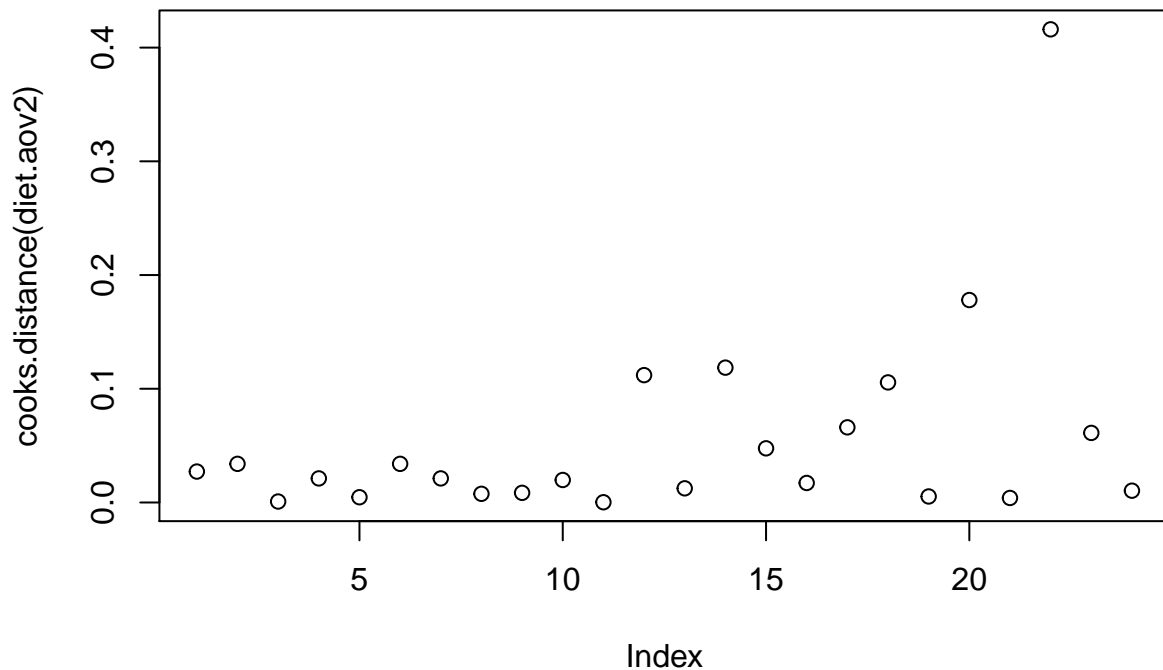


The residuals seem to be normally distributed

#### 7.4.3 Influential observations

We check the presence of influential observations by using Cook's distance.

```
plot(cooks.distance(diet.aov2))
```



```
diet_df[cooks.distance(diet.aov2)>0.3,]
```

```
##      LOSS JOGGING DIET
## 22    27      2km  fat
```

Observation number 22 has a large Cook's distance compared to the rest. We'll repeat the analysis for the data while deleting that point.

```
diet_df_small <- diet_df[-22,] # Delete observation 22
```

```
loss_small <- diet_df_small$LOSS
jogging_small <- diet_df_small$JOGGING
diet_small <- diet_df_small$DIET
```

Two-way ANOVA

```
diet.aov1_small <- aov(loss_small ~ jogging_small + diet_small + jogging_small * diet_small,
                       contrasts = list(jogging_small = "contr.sum", diet_small = "contr.sum"))
summary(diet.aov1_small)
```

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## jogging_small    2  542.0   271.00   70.977 5.16e-07 ***
## diet_small       3  204.3    68.09   17.832 0.000156 ***
## jogging_small:diet_small  6   15.5     2.58    0.675 0.672896
## Residuals      11   42.0     3.82
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Two-way ANOVA without interaction



```
diet.aov2_small <- aov(loss_small ~ jogging_small + diet_small,
                      contrasts = list(jogging_small = "contr.sum", diet_small = "contr.sum"))
summary(diet.aov2_small)
```

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## jogging_small  2  542.0   271.00    80.17 2.21e-09 ***
## diet_small     3  204.3    68.09    20.14 7.79e-06 ***
## Residuals     17   57.5     3.38
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

We obtain similar results as compared to before removing observation 22. We keep observation 22 in the data and continue the analysis.

## 7.5 Multiple comparisons for the main effects (in case interaction is not significant)

In our example, the interaction effect is not significant. The two main effects *JOGGING* and *DIET* are significant. We are going to use multiple comparisons techniques (Tukey, Scheffé, Bonferroni and Holm) on the main effects.

*Multiple comparisons on the main effects:*

```
library(agricolae)
library(sp)
diet.aov2 <- aov(LOSS ~ JOGGING + DIET, data = diet_df,
                contrasts = list(JOGGING = "contr.sum", DIET = "contr.sum"))
```

- Tukey HSD test for *JOGGING*

```
out <- HSD.test(diet.aov2, "JOGGING", group = FALSE)
out$means
```

```
##      LOSS      std r  Min  Max   Q25  Q50   Q75
## 0km 12.125 3.622844 8   7.5 16.5   8.50 12.5 15.50
## 1km 17.000 4.208834 8  11.0 23.0  13.75 17.0 20.25
## 2km 24.500 2.751623 8  19.5 27.5  23.50 24.5 27.00
```

```
out$comparison
```

```
##      difference pvalue signif.      LCL      UCL
## 0km - 1km      -4.875 0.0016    ** -7.858847 -1.891153
## 0km - 2km     -12.375 0.0000    *** -15.358847 -9.391153
## 1km - 2km      -7.500 0.0000    *** -10.483847 -4.516153
```

- Tukey HSD test for *DIET*

```
out1 <- HSD.test(diet.aov2, "DIET", group = FALSE)
out1$means
```

```
##      LOSS      std r  Min  Max   Q25  Q50   Q75
## carbo 20.00000 5.366563 6  13.5 27.5  16.125 19.5 23.625
## fat   14.83333 7.890923 6   7.5 27.0   9.125 12.0 19.750
## normal 15.66667 5.732946 6   8.5 24.5  12.125 15.0 18.625
## protein 21.00000 4.483302 6  15.5 27.0  17.375 21.5 23.750
```

```
out1$comparison
```

```
##      difference pvalue signif.      LCL      UCL
```

```
## carbo - fat      5.1666667 0.0062      **  1.3511428  8.982190
## carbo - normal   4.3333333 0.0229      *   0.5178095  8.148857
## carbo - protein -1.0000000 0.8794      -4.8155238  2.815524
## fat - normal     -0.8333333 0.9252      -4.6488572  2.982190
## fat - protein    -6.1666667 0.0012      ** -9.9821905 -2.351143
## normal - protein -5.3333333 0.0047      ** -9.1488572 -1.517810
```

### Remark:

- The majority of the R manuals suggest the Tukey HSD to ask for multiple comparisons test. But when writing a report, try to be systematic by using the same method for multiple comparisons everywhere.
- If we have a model with a significant interaction, before using the `HSD.test()` function, you first need to create a new variable that is equal to interaction. Then, apply an additive ANOVA model on three variables and after that apply the Tukey HSD test (or another test). (This will be discussed later).

## 7.6 Two-way ANOVA when cells have unequal sample size

### 7.6.1 What is an unbalanced design?

If you have a two-way ANOVA with **unequal numbers of entries per cell**, then the main effects and the interaction effect are no longer independent of each other. This type of design where the sample sizes for the different treatment combinations are not all equal is called an “**unbalanced design**”.

- With a *balanced design*, you have the following decomposition of the Total Sum of Squares:  

$$SSTO = SSA + SSB + SSAB + SSE$$
- In an *unbalanced design* this equation does not hold anymore.

Hence, the general recommendation for an unbalanced design is to use the regression approach:

- Use *Type III SS* to check the significance of the effects of the model.
- For the interpretation, use the least square averages instead of the sample averages.

### 7.6.2 Illustrative example

#### Example *Training*

The data *training.txt* contains information about children that were assigned to different training methods (*Method*) and that were separated for some period of time (*Sep\_Period*). The results of the test is registered in *score*.

	Length of separation period		
Method	20 minutes	40 minutes	60 minutes
No Training	26		6
	23	30	11
	28	25	17
	19	27	10
	18	36	14
			19
Training	15	24	31
	24	29	29
	25	23	35
	16	26	38
	22	27	34
	21	21	30

Import the data set *training.txt* as *training* in R.

```
training <- read.table(file=file.choose(), header=TRUE)
```

```
head(training)
```

```
##           Method Sep_Period score
## 1 No_Training    20_min    26
## 2 No_Training    20_min    23
## 3 No_Training    20_min    28
## 4 No_Training    20_min    19
## 5 No_Training    20_min    18
## 6 No_Training    40_min    30
```

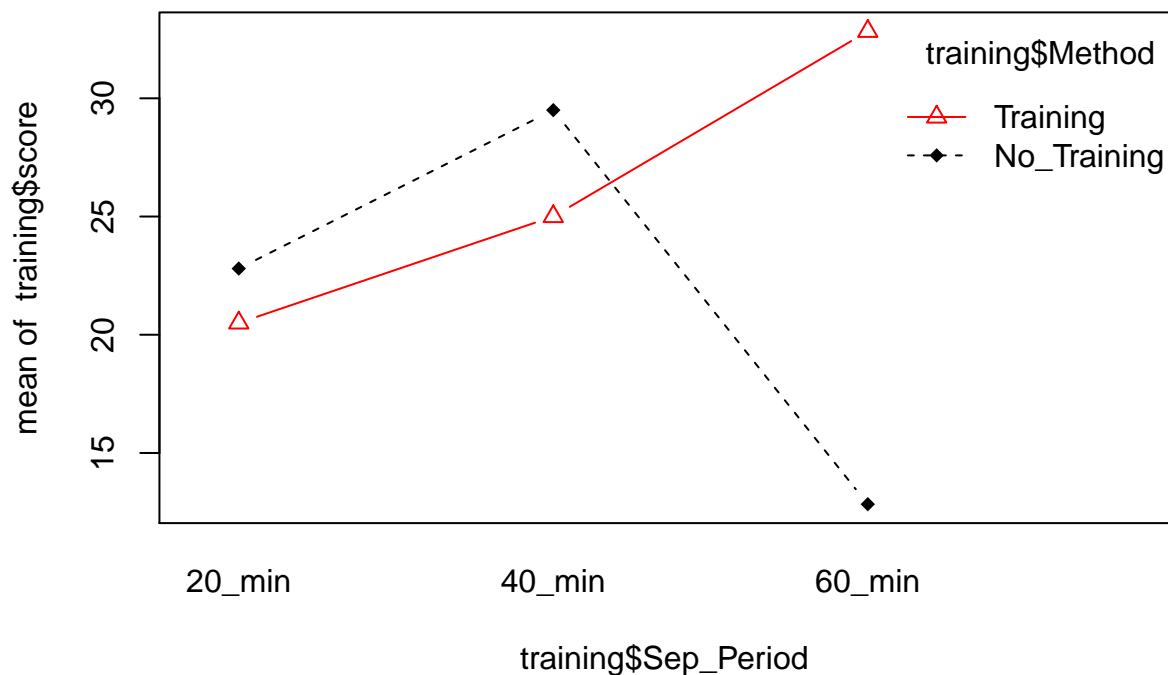
Descriptive statistics:

```
by_Method_SepPeriod <- group_by(training, Method, Sep_Period)
summarise(by_Method_SepPeriod, Avg = mean(score), SD = sd(score), number = n())
```

```
## # A tibble: 6 x 5
## # Groups:   Method [2]
##   Method      Sep_Period   Avg   SD number
##   <fct>       <fct>     <dbl> <dbl> <int>
## 1 No_Training 20_min      22.8  4.32     5
## 2 No_Training 40_min      29.5  4.80     4
## 3 No_Training 60_min      12.8  4.79     6
## 4 Training    20_min      20.5  4.14     6
## 5 Training    40_min      25    2.90     6
## 6 Training    60_min      32.8  3.43     6
```

Visualization:

```
library(stats)
interaction.plot(training$Sep_Period, training$Method, training$score, type="b",
                pch = c(18, 24, 22), col = c(1, 2, 3))
```



### 7.6.3 ANOVA table

For an unbalanced design, we use a regression approach. To obtain a *Type III SS ANOVA*, we have to do the following:

We are going to apply function `Anova()` from the `library(car)` in combination with function `lm()` for the linear models.

- Set for `lm()` the contrast from `contr.treatment` to `contr.sum`.

- Specify for Anova() the value of type = "III".

```
# ANOVA table in case of unbalanced design: use of lm function
training.lm <- lm(score ~ Method + Sep_Period + Method*Sep_Period, data = training,
                  contrasts = list(Method = "contr.sum", Sep_Period = "contr.sum"))
Anova(training.lm, type = "III")
```

```
## Anova Table (Type III tests)
##
## Response: score
##              Sum Sq Df    F value    Pr(>F)
## (Intercept)  18432.3  1 1118.4453 < 2.2e-16 ***
## Method        156.0  1   9.4681  0.004753 **
## Sep_Period    175.8  2   5.3333  0.011168 *
## Method:Sep_Period 1036.3  2  31.4404  8.89e-08 ***
## Residuals      445.0 27
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

**Step 1:** Check whether the interaction is significant

We can see that the interaction is significant ( $p\text{-value} < 0.05$ ). It means that we are not allowed to interpret the main effects.

**Step 2A:**

- Check the diagnostics.
- Use pairwise comparisons on interaction effect.

Both will be treated in the following two sections.

#### 7.6.4 Diagnostics

1. Check assumption of *normality*

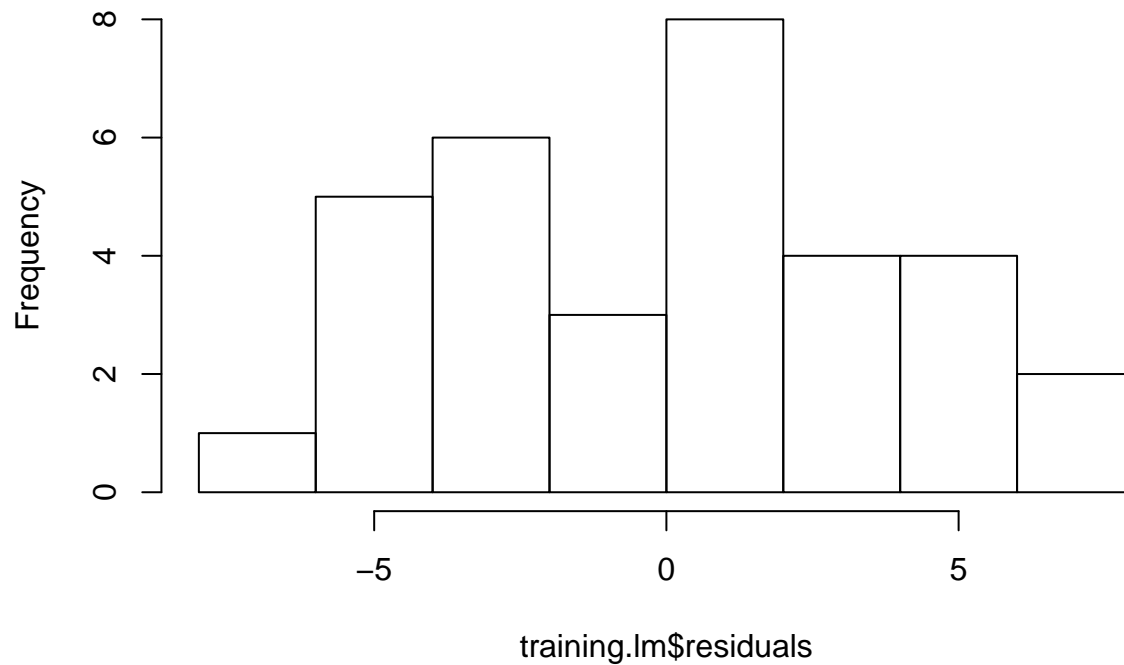
Test normality of the within-cell residuals

```
shapiro.test(training.lm$residuals)
```

```
##
## Shapiro-Wilk normality test
##
## data:  training.lm$residuals
## W = 0.96223, p-value = 0.2985
```

```
hist(training.lm$residuals)
```

## Histogram of training.lm\$residuals



2. Check assumption of *homogeneity of variances*

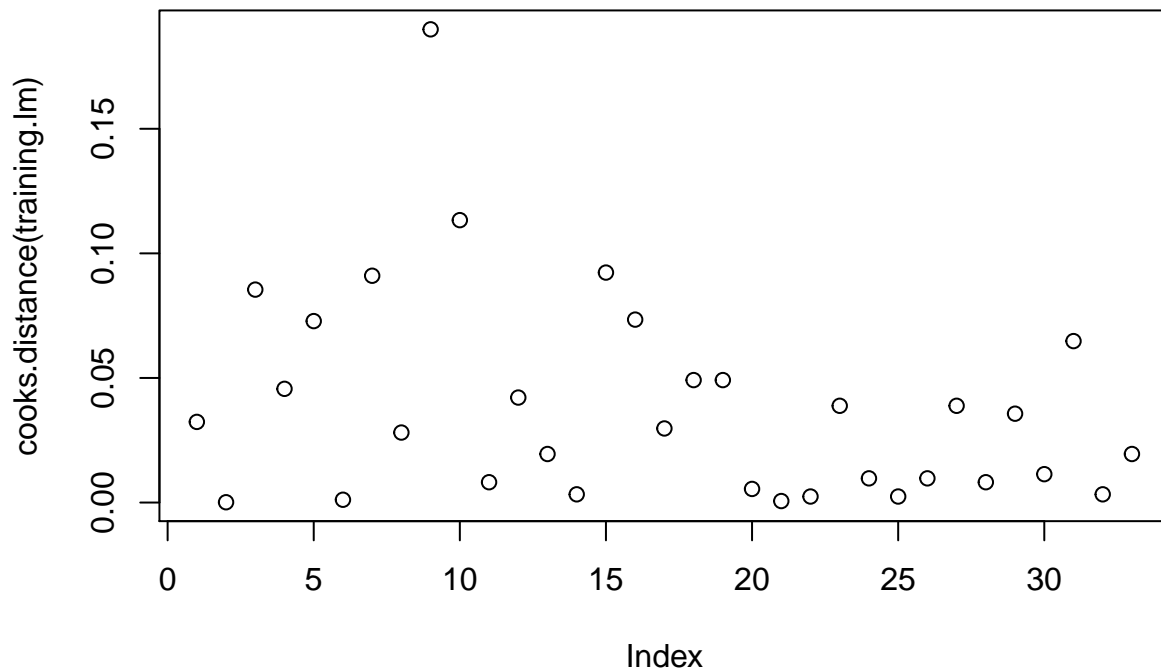
```
leveneTest(score ~ Method*Sep_Period, data = training)
```

```
## Levene's Test for Homogeneity of Variance (center = median)
##      Df F value Pr(>F)
## group  5  0.3678 0.8661
##      27
```

3. Check *influential observations*

Plotting Cook's distance

```
plot(cooks.distance(training.lm))
```



### 7.6.5 Pairwise comparisons of treatment effects

Since the interaction term is significant, we are interested in the pairwise comparisons of the interaction effect (*Method \* Sep\_Period*).

#### **Method:**

Create one new variable which is the interaction term. Refit a one-way ANOVA with as single variable this interaction term. Then you can use Tukey and other MC (multiple comparison) methods on this new variable.

```
## Method 1: Create a new variable with the interaction term
# Initialize the new variable
method_sep_period <- character(length(training$score))
# Store values in new variable
for (i in 1:length(training$score))
{method_sep_period[i] <- paste(substring(training$Method[i], 1, 4),
                              training$Sep_Period[i], sep="")}
# Create a new data frame
new_df <- data.frame(score = training$score, method_sep_period)
head(new_df, n = 10)
```

```
##   score method_sep_period
## 1    26      No_T20_min
## 2    23      No_T20_min
## 3    28      No_T20_min
## 4    19      No_T20_min
## 5    18      No_T20_min
## 6    30      No_T40_min
```

```
## 7      25      No_T40_min
## 8      27      No_T40_min
## 9      36      No_T40_min
## 10     6      No_T60_min
```

Descriptive statistics

```
describe <- describeBy(new_df$score, new_df$method_sep_period, mat = TRUE)
describe.st <- subset(describe, select=c("group1", "n", "mean", "sd", "median", "min", "max"))
describe.st
```

```
##      group1 n      mean      sd median min max
## X11 No_T20_min 5 22.80000 4.324350  23.0  18  28
## X12 No_T40_min 4 29.50000 4.795832  28.5  25  36
## X13 No_T60_min 6 12.83333 4.792355  12.5   6  19
## X14 Trai20_min 6 20.50000 4.135215  21.5  15  25
## X15 Trai40_min 6 25.00000 2.898275  25.0  21  29
## X16 Trai60_min 6 32.83333 3.430258  32.5  29  38
```

```
# Apply one-way ANOVA on this new data frame
new_df.aov <- aov(score ~ method_sep_period, new_df,
                  contrasts = list(method_sep_period = "contr.sum"))
summary(new_df.aov)
```

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## method_sep_period  5    1419   283.78    17.22 1.17e-07 ***
## Residuals        27     445    16.48
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# Ask for Tukey HSD test
out2 <- HSD.test(new_df.aov, "method_sep_period", group = FALSE)
round(out2$means, 2)
```

```
##      score std r Min Max  Q25  Q50  Q75
## No_T20_min 22.80 4.32 5  18  28 19.00 23.0 26.00
## No_T40_min 29.50 4.80 4  25  36 26.50 28.5 31.50
## No_T60_min 12.83 4.79 6   6  19 10.25 12.5 16.25
## Trai20_min 20.50 4.14 6  15  25 17.25 21.5 23.50
## Trai40_min 25.00 2.90 6  21  29 23.25 25.0 26.75
## Trai60_min 32.83 3.43 6  29  38 30.25 32.5 34.75
```

```
out2$comparison
```

```
##      difference pvalue signif.      LCL      UCL
## No_T20_min - No_T40_min -6.700000 0.1718      -15.0436290  1.6436290
## No_T20_min - No_T60_min  9.966667 0.0046      **  2.4351153 17.4982181
## No_T20_min - Trai20_min  2.300000 0.9336      -5.2315514  9.8315514
## No_T20_min - Trai40_min -2.200000 0.9444      -9.7315514  5.3315514
## No_T20_min - Trai60_min -10.033333 0.0043      ** -17.5648847 -2.5017819
## No_T40_min - No_T60_min 16.666667 0.0000      ***  8.6380059 24.6953274
## No_T40_min - Trai20_min  9.000000 0.0213      *  0.9713392 17.0286608
## No_T40_min - Trai40_min  4.500000 0.5328      -3.5286608 12.5286608
## No_T40_min - Trai60_min -3.333333 0.7972      -11.3619941  4.6953274
## No_T60_min - Trai20_min -7.666667 0.0313      * -14.8477191 -0.4856142
## No_T60_min - Trai40_min -12.166667 0.0002      *** -19.3477191 -4.9856142
## No_T60_min - Trai60_min -20.000000 0.0000      *** -27.1810525 -12.8189475
## Trai20_min - Trai40_min -4.500000 0.4123      -11.6810525  2.6810525
```



```
## Trai20_min - Trai60_min -12.333333 0.0002      *** -19.5143858 -5.1522809
## Trai40_min - Trai60_min -7.833333 0.0265      * -15.0143858 -0.6522809
```

## 8 Experimental design

### 8.1 Observational study versus designed experiment

- In an **experiment** investigators apply treatments to experimental units (people, animals, plots of land, etc.) and then proceed to observe the effect of the treatments on the experimental units.
- In an **observational study**, investigators observe subjects and measure variables of interest without assigning treatments to the subjects. The treatment that each subject receives is determined beyond the control of the investigator.

#### Example *Smoking*

Suppose we want to study the effect of smoking on the lung capacity in women.

##### Experiment

- Find 100 women age 30 who do not currently smoke.
- Randomly assign 50 of the 100 women to the smoking treatment and the other 50 to the no-smoking treatment.
- Those in the smoking group smoke a pack a day for 10 years while those in the control group remain smoke free for 10 years.
- Measure lung capacity for each of the 100 women.
- Analyze, interpret and draw conclusions from data.

##### Observational study

- Find 100 women age 40 of which 50 have been smoking a pack a day for 10 years while the other 50 have been smoke free for 10 years.
- Measure lung capacity for each of the 100 women.
- Analyze, interpret and draw conclusions from data.

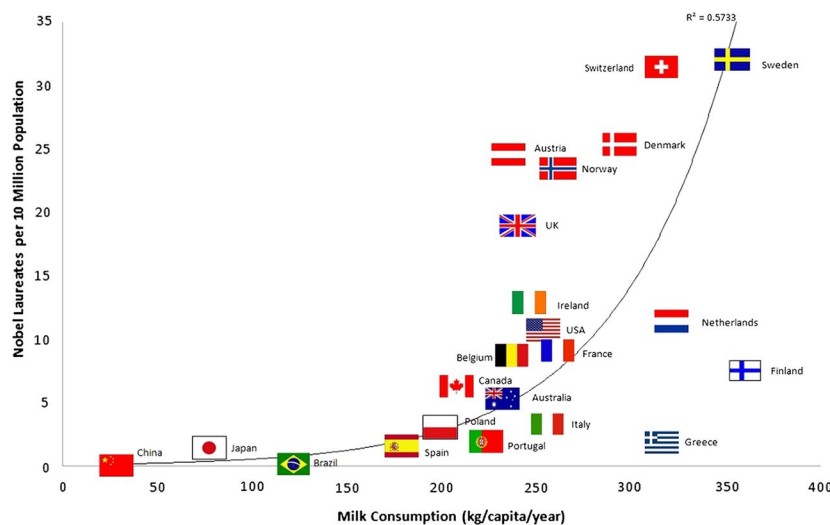
##### Fisher's Hypothesis

- Suppose there is a gene that causes smoking to appear to be a very pleasurable experience.
- Suppose the same gene also causes emphysema, lung cancer, throat cancer, etc.
- People who have the gene will be more likely to smoke than people who do not have the gene.
- People who have the gene will be more likely to get emphysema, lung cancer, throat cancer, etc
- So is it really smoking that causes health problems? Maybe it is just the gene?
- A **confounding variable** is related both to group membership and to the outcome of interest. Its presence makes it hard to establish the outcome as being a direct consequence of group membership.

Correlation does not imply causation!



Correlation between countries' annual milk consumption (kg/capita/year) and the number of Nobel laureates per 10 million population.



## 8.2 Basic principles of experimental design

### 8.2.1 Replication

In all experiments, some variation is introduced because of the fact that the experimental units such as individuals or plots of land in agricultural experiments cannot be physically identical. This type of variation

can be removed by using a number of experimental units. We therefore perform the experiment more than once, i.e., we repeat the basic experiment.

**Replication allows us to estimate the experimental error** and to perform statistical analysis.

### 8.2.2 Randomization

Randomization is a random process of assigning treatments to the experimental units.

**The purpose of randomization is to remove bias and other sources of uncontrollable variation.**

Another advantage of randomization (accompanied by replication) is that it forms the basis of any valid statistical test. Hence the treatments must be assigned at random to the experimental units.

#### Example *Corn yield*

You want to compare the yield for two types of corn (type *A* and type *B*).

We have several small fields which are available, but the fertility at one side of the land is different from the fertility at the other side.

First suggestion of assigning the different types to the 10 subfields:

<b>A</b>	<b>A</b>	<b>B</b>	<b>B</b>	<b>B</b>
<b>A</b>	<b>A</b>	<b>A</b>	<b>B</b>	<b>B</b>

Problem here: If we detect a difference in the yield, we cannot detect whether it comes from the different type of corn or whether it comes from the difference in fertility of the ground.

**Remark:** systematic arranging the type of corn over the several plots  $\neq$  randomization

<b>A</b>	<b>B</b>	<b>A</b>	<b>B</b>	<b>A</b>
<b>B</b>	<b>A</b>	<b>B</b>	<b>A</b>	<b>B</b>

### 8.2.3 Blocking

It has been observed that all sources of uncontrollable variation are not removed by randomization and replication. A block is a subset of experimental conditions that are expected to be more homogeneous than the rest.


Blocking refers to the method of creating homogeneous blocks of data in which the nuisance factor is kept constant and the factor of interest is allowed to vary.

**Blocking is used to eliminate the variability due to the difference between block.**

#### Example *Corn yield*

Example of blocking

Field 1	Field 2	Field 3	Field 4	Field5
A	B	A	A	B
B	A	B	B	A



**Increasing fertility of the ground**

Within each block, the types are randomly assigned.

**Remark:**

Both blocking and randomization deal with nuisance <sup>1</sup> factors.

- Blocking can only be used when the nuisance factor is under our control (e.g. choice of materials or substances).
- If the nuisance factor is not under our control, then randomization remains the only tool available.

‘Block what you can, randomize what you cannot!’

## 9 The general linear model

In previous chapters we have seen

- Linear regression (simple and multiple)
- ANOVA (one-way and two-way)

The above models are special cases of the **General Linear Model** (GLM).

Models with continuous response variable

Explanatory variables	Response variable	Method
Continuous	Continuous	Regression
Categorical	Continuous	ANOVA
Continuous and categorical	Continuous	GLM

<sup>1</sup>A nuisance factor is a factor that has some effect on the response, but is of no interest to the experimenter. However, the variability it transmits to the response needs to be minimized or explained. Hence, nuisance factors needs to be taken into account in an analysis.