

# MA50259: Statistical Design of Investigations

## *Lab sheet 1: Completely randomised design*

In this practical you will learn how to simulate and analyse your own completely randomised experimental design. Take the time to run each of the following commands and analyse the displayed results to understand what the code is doing.

1. Install and load the `tidyverse` package

```
install.packages('tidyverse')
```

```
library('tidyverse')
```

2. Set the number of treatment levels, `t` to 6 and the number of replicates per treatment level `r` (balanced design) to 10. This determines the number of experimental units `n`, that is, 60.

```
r<-10  
t<-6  
n<-t*r
```

3. Create a character vector with the names of the treatment levels

```
levels<-c("level 1","level 2","level 3","level 4", "level 5","level 6")  
levels
```

```
## [1] "level 1" "level 2" "level 3" "level 4" "level 5" "level 6"
```

```
class(levels)
```

```
## [1] "character"
```

4. Create a vector of length `n` with `r` replicates of each of the `t` treatment levels

```
v<-rep(levels,each = r)  
v
```

```
## [1] "level 1" "level 1" "level 1" "level 1" "level 1" "level 1" "level 1" "level 1"  
## [8] "level 1" "level 1" "level 1" "level 2" "level 2" "level 2" "level 2" "level 2"  
## [15] "level 2" "level 2" "level 2" "level 2" "level 2" "level 2" "level 2" "level 3"  
## [22] "level 3" "level 3" "level 3" "level 3" "level 3" "level 3" "level 3" "level 3"  
## [29] "level 3" "level 3" "level 4" "level 4" "level 4" "level 4" "level 4" "level 4"  
## [36] "level 4" "level 4" "level 4" "level 4" "level 4" "level 4" "level 5" "level 5"  
## [43] "level 5" "level 5" "level 5" "level 5" "level 5" "level 5" "level 5" "level 5"  
## [50] "level 5" "level 6" "level 6" "level 6" "level 6" "level 6" "level 6" "level 6"  
## [57] "level 6" "level 6" "level 6" "level 6" "level 6" "level 6" "level 6" "level 6"
```

```
class(v)
```

```
## [1] "character"
```

5. Convert the vector `v` into a factor that R can manage for statistical analysis

```
f <- factor(v)  
f
```

```
## [1] level 1 level 1 level 1 level 1 level 1 level 1 level 1 level 1 level 1  
## [9] level 1 level 1 level 2 level 2 level 2 level 2 level 2 level 2 level 2  
## [17] level 2 level 2 level 2 level 2 level 3 level 3 level 3 level 3 level 3  
## [25] level 3 level 3 level 3 level 3 level 3 level 3 level 4 level 4 level 4
```

```
## [33] level 4 level 4 level 4 level 4 level 4 level 4 level 4 level 4
## [41] level 5 level 5 level 5 level 5 level 5 level 5 level 5 level 5
## [49] level 5 level 5 level 6 level 6 level 6 level 6 level 6 level 6
## [57] level 6 level 6 level 6 level 6
## Levels: level 1 level 2 level 3 level 4 level 5 level 6
```

```
class(f)
```

```
## [1] "factor"
```

6. Two alternative ways to create the vector `f` are as follows

```
# alternative code 1
f <- factor(rep(levels,each = r))
# Nesting the 2 commands avoids the unnecessary intermediate vector v

# alternative code 2
f <- rep(levels,each = r) %>% factor()
# Pipe %>% is equivalent the nested commands above but makes clear what is evaluated first
```

7. The following commands perform the randomisation of the experimental units to treatment levels as a random re-shuffling of the factor `f`. You may want to have a look at the help file of the R command `sample`

```
set.seed(5678)
1:n # row labels
```

```
## [1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23
## [24] 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46
## [47] 47 48 49 50 51 52 53 54 55 56 57 58 59 60
```

```
rows<-sample(1:n,n) # random re-shuffling of rows
rows
```

```
## [1] 24 26 60 22 58 6 10 19 8 46 4 49 57 7 45 32 28 15 40 9 47 14 50
## [24] 41 37 2 23 35 1 3 54 55 34 59 39 51 43 30 31 27 17 33 20 56 29 5
## [47] 21 38 12 52 44 16 13 11 18 36 42 53 48 25
```

```
fac <- f[rows]
fac
```

```
## [1] level 3 level 3 level 6 level 3 level 6 level 1 level 1 level 2
## [9] level 1 level 5 level 1 level 5 level 6 level 1 level 5 level 4
## [17] level 3 level 2 level 4 level 1 level 5 level 2 level 5 level 5
## [25] level 4 level 1 level 3 level 4 level 1 level 1 level 6 level 6
## [33] level 4 level 6 level 4 level 6 level 5 level 3 level 4 level 3
## [41] level 2 level 4 level 2 level 6 level 3 level 1 level 3 level 4
## [49] level 2 level 6 level 5 level 2 level 2 level 2 level 2 level 4
## [57] level 5 level 6 level 5 level 3
## Levels: level 1 level 2 level 3 level 4 level 5 level 6
```

```
# set.seed is used to control randomness if we do not set the seed then we obtain different results
# everytime we run the program!!! You can test that yourself!
```

8. An alternative code for creating the factor `fac` is as follows (as it is presented in the lecture slides):

```
set.seed(5678)
fac <- sample(f,size=n)
# fac <- f %>% sample(size=n) alternative code
```

9. Now, create a data.frame (a dataset, basically) with two variables: the (randomised) factor `fac` and a vector of simple labels for the units, namely, from 1 to n

```
units <- 1:n # unit labels
crd <- tibble(units=units, treatment=fac ) # creates data.frame (called tibble in the tidyverse)
crd
```

```
## # A tibble: 60 x 2
##   units treatment
##   <int> <fct>
## 1     1 level 3
## 2     2 level 3
## 3     3 level 6
## 4     4 level 3
## 5     5 level 6
## 6     6 level 1
## 7     7 level 1
## 8     8 level 2
## 9     9 level 1
## 10    10 level 5
## # ... with 50 more rows
```

```
glimpse(crd)
```

```
## Observations: 60
## Variables: 2
## $ units      <int> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 1...
## $ treatment <fct> level 3, level 3, level 6, level 3, level 6, level 1...
```

10. Re-arrange the rows of the data.frame `crd` so that experimental units with the same treatment levels are together and the treatments levels are in their original order: level 1, level 2, etc.

```
crd<-arrange(crd,treatment) # arrange by treatment level
# alternative code
# crd<- crd %>% arrange(treatment)
glimpse(crd)
```

```
## Observations: 60
## Variables: 2
## $ units      <int> 6, 7, 9, 11, 14, 20, 26, 29, 30, 46, 8, 18, 22, 41, ...
## $ treatment <fct> level 1, level 1, level 1, level 1, level 1, level 1...
```

11. Note we did the labeling of experimental units after the randomisation! This might not be appropriate if (for any reason) we need to identify some experimental units with their original labels. The following code performs this and creates a new data.frame called `crd1`

```
crd1 <- tibble(units=units, treatment=f )
crd1<-crd1[rows,]
crd1
```

```
## # A tibble: 60 x 2
##   units treatment
##   <int> <fct>
## 1    24 level 3
## 2    26 level 3
## 3    60 level 6
## 4    22 level 3
## 5    58 level 6
```

```
## 6      6 level 1
## 7     10 level 1
## 8     19 level 2
## 9      8 level 1
## 10    46 level 5
## # ... with 50 more rows

crd1<- crd1 %>% arrange(treatment)
crd1
```

```
## # A tibble: 60 x 2
##   units treatment
##   <int> <fct>
## 1      6 level 1
## 2     10 level 1
## 3      8 level 1
## 4      4 level 1
## 5      7 level 1
## 6      9 level 1
## 7      2 level 1
## 8      1 level 1
## 9      3 level 1
## 10     5 level 1
## # ... with 50 more rows
```

Note the two data.frames `crd` and `crd1` are different since the experimental units are labelled differently.

12. Now we simulate the response values using a simple model. We set the mean response for each treatment levels as follows: first set the overall mean response  $\mu$  and also set the effects  $\tau_i$  of each treatment level. The mean response values for each treatment level is given by  $\mu + \tau_i$

```
mu<-500 # overall mean response
tau<-c(-20,50,0,-30,-10,100) # treatment effects over mu
means<-mu+tau # mean response for each treatment level
means
```

```
## [1] 480 550 500 470 490 600
```

13. Create a vector of length `n` with `r` replicates of each of the `t` different mean response values

```
means<-mu+tau %>% rep(each=r) # vector of means
# means<-rep(mu+tau,each=r) # alternative code
means
```

```
## [1] 480 480 480 480 480 480 480 480 480 480 550 550 550 550 550 550 550
## [18] 550 550 550 500 500 500 500 500 500 500 500 500 500 470 470 470 470
## [35] 470 470 470 470 470 470 490 490 490 490 490 490 490 490 490 490 600
## [52] 600 600 600 600 600 600 600 600 600
```

14. Set the standard deviation  $\sigma$  and simulate the response values for all the experimental units by simulating Gaussian (Normal) random numbers with means given by `means` and common standard deviation  $\sigma$

```
sd<-10 # common standard deviation
y<-rnorm(n,mean=means,sd=sd) # note the units should be arranged by treatment level
y
```

```
## [1] 472.4084 485.5443 469.4177 476.5448 465.7961 478.7417 467.5876
## [8] 493.3207 478.1360 475.6323 556.9446 529.1966 547.6057 544.2420
## [15] 549.0948 561.5796 560.7574 561.1673 551.7749 554.6377 492.2205
## [22] 493.1161 498.0180 496.2580 485.5763 508.7216 504.7826 501.5454
```

```
## [29] 515.4524 495.0196 461.8719 458.5625 472.9068 459.4350 456.0150
## [36] 480.9937 465.8724 468.7799 456.7203 470.4136 489.8517 465.4183
## [43] 468.4639 484.2266 495.4763 485.2693 485.9455 473.3897 471.6645
## [50] 484.7507 599.8706 591.0384 618.1664 607.3093 614.9173 581.2274
## [57] 587.9494 596.5072 599.7507 586.7642
```

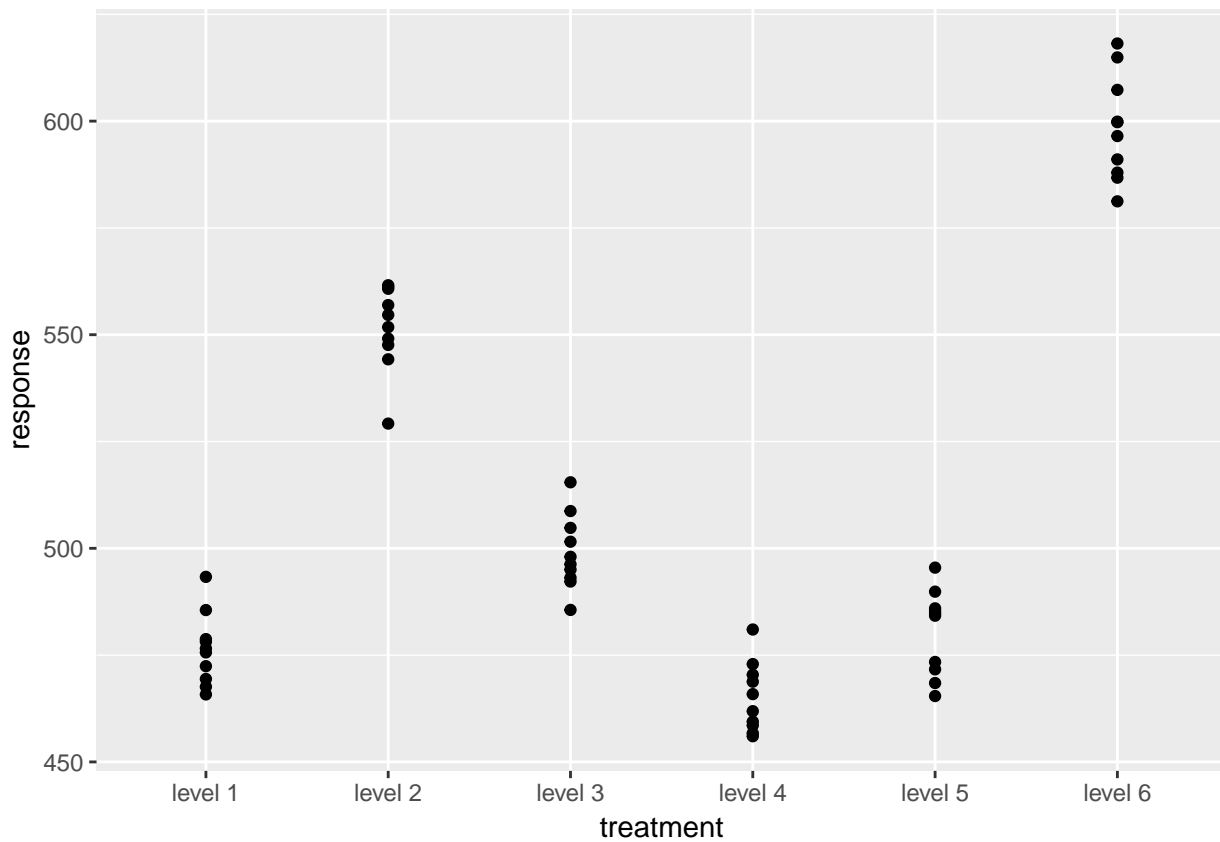
15. Append the generated vector of response values to the `crd1` data.frame

```
crd1$response<-y
glimpse(crd1)
```

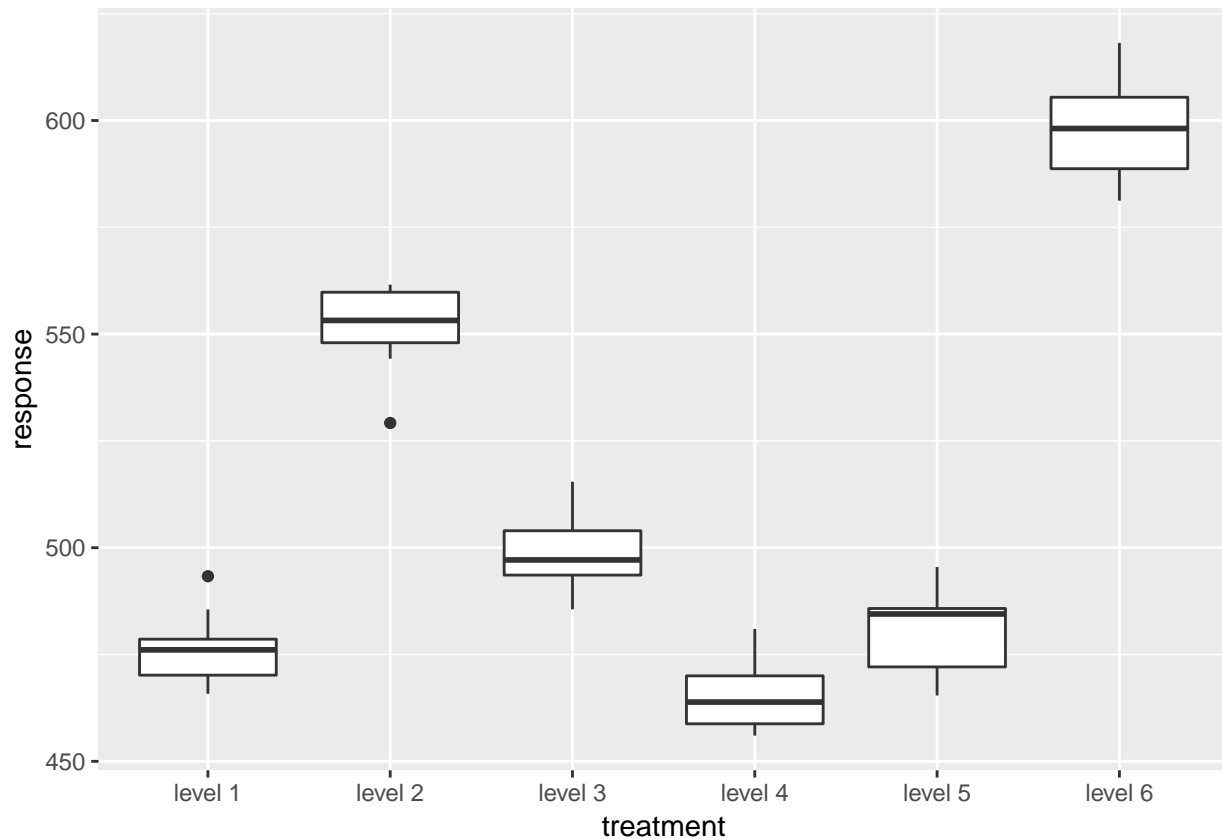
```
## Observations: 60
## Variables: 3
## $ units      <int> 6, 10, 8, 4, 7, 9, 2, 1, 3, 5, 19, 15, 14, 17, 20, 1...
## $ treatment  <fct> level 1, level 1, level 1, level 1, level 1, level 1...
## $ response   <dbl> 472.4084, 485.5443, 469.4177, 476.5448, 465.7961, 47...
```

16. Plot the data in `crd1` using `ggplot`

```
ggplot(crd1,aes(treatment,response))+geom_point()
```



```
ggplot(crd1,aes(treatment,response))+geom_boxplot()
```



17. Analyse the data using a linear model via the `lm` command in R

```
mod.crd1<-lm(response~treatment,data=crd1)
coefs<-coef(mod.crd1)
coefs
```

```
##      (Intercept) treatmentlevel 2 treatmentlevel 3 treatmentlevel 4
##      476.312968      75.387111      22.758085      -11.155876
## treatmentlevel 5 treatmentlevel 6
##      4.132686      122.037112
```

The estimates given, correspond to estimates of  $\mu$ ,  $\tau_2$ ,  $\tau_3$ ,  $\tau_4$ ,  $\tau_5$  and  $\tau_6$  so you can compare them with the actual population values.

```
c(mu,tau[-1])
```

```
## [1] 500  50   0 -30 -10 100
```

The estimate of  $\tau_1$  has been arbitrarily set to zero!

18. The estimates of the mean values for each level of the treatment effect are given as follows

```
taus<-c(0,coefs[-1])
means2<-coefs[1]+taus
means2
```

```
##      treatmentlevel 2 treatmentlevel 3 treatmentlevel 4
##      476.3130      551.7001      499.0711      465.1571
## treatmentlevel 5 treatmentlevel 6
##      480.4457      598.3501
```

You can compare these with the actual population values

```
mu+tau
```

```
## [1] 480 550 500 470 490 600
```

19. It turns out that the above mean estimates are exactly the same as the sample response means for each level of the treatment factor.

```
by_group <- group_by(crd1, treatment) #groups the data frame  
by_group
```

```
## # A tibble: 60 x 3  
## # Groups:   treatment [6]  
##   units treatment response  
##   <int> <fct>      <dbl>  
## 1     6 level 1      472.  
## 2    10 level 1      486.  
## 3     8 level 1      469.  
## 4     4 level 1      477.  
## 5     7 level 1      466.  
## 6     9 level 1      479.  
## 7     2 level 1      468.  
## 8     1 level 1      493.  
## 9     3 level 1      478.  
## 10    5 level 1      476.  
## # ... with 50 more rows
```

```
summaries.crd<-summarize(by_group, means = mean(response))  
# applies the mean function for each group defined by the levels of the treatment factor  
glimpse(summaries.crd)
```

```
## Observations: 6  
## Variables: 2  
## $ treatment <fct> level 1, level 2, level 3, level 4, level 5, level 6  
## $ means      <dbl> 476.3130, 551.7001, 499.0711, 465.1571, 480.4457, 59...
```