

# Computer exercises on the BOOTSTRAP

Exercise 1:

Create a dataframe which contains information on normal distributed observations measured in groups 1 ( $n_1=23$ ) and 2 ( $n_2=25$ ).

```
my.data<-  
data.frame(group=c(rep(1,23),rep(2,25)),values=c(rnorm(23,1,1),rnorm(25,2.5,1)))
```

Give some descriptive statistics of the observed data:

```
group.1.summary<-summary(my.data$value[my.data$group==1])
```

What kind of object is `group.1.summary` ? Calculate a `group.1.summary`. Make a table which contains the summaries of the observed data.

```
describe.stat.table<-rbind(group.1.summary,group.2.summary)
```

Give some nice names to the Row-names of the object `describe.stat.table`.

```
> dimnames(describe.stat.table)[[1]]<-c("Group 1","Group 2")  
> describe.stat.table  
      Min.   1st Qu.  Median    Mean 3rd Qu.   Max.  
Group 1 -0.8754 -0.003584 0.5653 0.7326   1.441 2.483  
Group 2  0.3748  1.514000 2.1910 2.1380   2.507 4.333
```

The standard deviation of the measurement is not given in the table. Please add it to the information (with the same amount of decimals as the other numbers).

```
STD<-  
sqrt(c(var(my.data$value[my.data$group==1]),var(my.data$value[my.data$group==2])))  
describe.stat.table<-cbind(describe.stat.table,round(STD,3))
```

## Exercise 2: Bootstrap t-test

Use the function `twosam` to calculate the t-statistics for the data in your dataframe `my.data`.

```
t.obs<-twosam(my.data$value[my.data$group==1], my.data$value[my.data$group==2])
```

Write the following function to imitate the unknown random process:

```
> boot.two.sample <- function (i,obs,gr)
  {
    l<-length(gr)
    ss<-sample(l,l,replace=T)
    gr.new<-gr[ss]
    x.1<-obs[gr.new==1]
    x.2<-obs[gr.new==2]
    return(twosam(x.1,x.2))
  }
```

Run a bootstrap with 999 samples:

```
boot.res<-unlist(lapply(1:999,boot.two.sample,obs=my.data$value,gr=my.data$group))
```

Perform a two-sided bootstrap t-test by calculating the p-value:

```
p.boot<-sum( ifelse(abs(boot.res)>abs(t.obs),1,0) )/999
```

Interpret the result.

Compare the bootstrap result to a standard t-test which can be calculated by

```
t.test(my.data$value[my.data$group==1],my.data$value[my.data$group==2])
```

### Exercise 3: Bootstrap confidence intervals

Read from the result of the standard t-test in R the estimate for the difference of the group means and its 95% confidence interval.

```
t.test.res<- t.test(my.data$value[my.data$group==1],my.data$value[my.data$group==2])
diff.obs<- t.test.res$estimate[1]- t.test.res$estimate[2]
```

Calculate a bootstrap sample of mean differences by using the function

```
> boot.mean <- function (i,obs,gr)
  {
    l<-length(gr)
    ss<-sample(l,l,replace=T)
    gr.new<-gr[ss]
    x.1<-obs[gr.new==1]
    x.2<-obs[gr.new==2]
    return(mean(x.1)-mean(x.2))
  }
```

Make a bootstrap sample of mean differences with 999 replicates.

```
mean.res<-unlist(lapply(1:999,boot.mean,obs=my.data$value,gr=my.data$group))
```

Use two approaches to calculate a 95% bootstrap confidence interval for the mean difference.

- 1.) calculate the variance of the bootstrap sample  $v^2$  [ $\text{var}(\text{mean.res})$ ] and use  
 $[\text{diff.obs} - 1.96 \cdot v; \text{diff.obs} + 1.96 \cdot v]$
- 2.) order the bootstrap sample and take the 25<sup>th</sup> element,  $\text{diff}_{0.025}$ , and the 975<sup>th</sup> element,  
 $\text{diff}_{0.975}$  to get  $[2 \cdot \text{diff.obs} - \text{diff}_{0.975}; 2 \cdot \text{diff.obs} - \text{diff}_{0.025}]$

```
mean.res<-sort(mean.res)
diff.975<- mean.res[975]
diff.025<- mean.res[25]
```

Compare the bootstrap confidence intervals with the exact 95% confidence intervals. Is the true mean difference of  $1 - 2.5 = -1.5$  included in the confidence intervals?

## Exercise 4: Balanced 2 by 2 ANOVA

Perform for the two.by.two data an ANOVA to study the Mutation by treatment effect.

The data looks as follows and gives VSN transformed values of signal intensities:

```
> two.by.two[1,]
  W.NT.1  W.TR.1  W.NT.2  W.TR.2  W.NT.3  W.TR.3  W.NT.4  W.TR.4
12.64218 12.92974 12.97218 13.28092 12.71829 13.25633 13.17850 12.89417
  W.NT.5  W.TR.5  M.NT.6  M.TR.6  M.NT.7  M.TR.7  M.NT.8  M.TR.8
13.33704 12.57546 13.10308 12.75008 12.57069 12.51956 12.83289 12.47562
  M.NT.9  M.TR.9  M.NT.10 M.TR.10
13.02821 12.71112 12.28821 12.37065
```

W – wild type / M – mutation

NT – not treated / TR – treated

There are 10 treated and 10 non-treated samples, as well as 10 wild type and 10 mutations.

The function `two.by.two.anova.rfc` performs the analysis and returns a list with two components: the residuals and the interesting `ge.wt.nt` values.

```
two.by.two.anova.res<-two.by.two.anova.rfc()
```

Plot a histogram of the residuals:

```
hist(two.by.two.anova.res[[1]],main="Histogram of residuals")
```

Calculate the standard deviation of the residuals. `sqrt(var(two.by.two.anova.res[[1]])`

Look at the differential gene expression with respect to treatment between wild type and mutations. What is the minimum, what the maximum log fold change?

The fold-change for a differential expression for gene *i* is calculated by

```
exp(2 * two.by.two.anova.res[[2]][i])
```

How many genes would be differentially expressed if one uses the 95% CI based on  $1.96 \cdot SE$ ?

```
table(2*abs(two.by.two.anova.res[[2]])>1.96*sqrt(var(two.by.two.anova.res$resid)))
```