

Report Template coursework assignment A - 2018

CS4125 Seminar Research Methodology for Data Science

Pórunn Arna Ómardóttir (4917499), Nathan Buskulic (4947916), Mitchell Deen(4396340)

4/3/2019

Contents

1	Part 1 - Design and set-up of true experiment	1
1.1	The motivation for the planned research.	1
1.2	The theory underlying the research.	1
1.3	Research questions	2
1.4	The related conceptual model	2
1.5	Experimental Design	2
1.6	Experimental procedure	2
1.7	Measures	2
1.8	Participants	2
1.9	Suggested statistical analyses	2
2	Part 2 - Generalized linear models	3
2.1	Question 1 Twitter sentiment analysis (Between groups - single factor)	3
2.1.1	Collecting tweets, and data preparation	3
2.1.2	Conceptual model	3
2.1.3	Homogeneity of variance analysis	3
2.1.4	Visual inspection	4
2.1.5	Mean sentiments	6
2.1.6	Linear model	7
2.1.7	Post Hoc analysis	7
2.1.8	Report section for a scientific publication	8

1 Part 1 - Design and set-up of true experiment

1.1 The motivation for the planned research.

(Max 250 words) The coffee is today the most consumed drink in the world and it is told to increase your performance and concentration. We want to challenge this idea and verify scientifically if this is a valid idea. We want to test how coffee consumption (and the level of caffeine inside) affect the result of an IQ test. We are most interesting and seeing what the affect is on TU Delft students like ourselves. So the participants will be recruited from the TU Delft student body. # Add that we are doing that on tudelft student

1.2 The theory underlying the research.

(Max 250 words) Preferable based on theories reported in literature There is a large body of literature available on the effects of caffeine on the performance in cognitive tasks. Literature generally supports the idea that coffee improves this performance, see e.g. (Jarvis, 1993; Nehlig, 2010; Rogers et al., 2008). In a brief survey of the relevant literature we did not find any studies specifically addressing students. We would like to investigate this part of the population in more detail.

Jarvis, M. J. (1993). Does caffeine intake enhance absolute levels of cognitive performance?. Psychopharmacology, 110(1-2), 45-52. Rogers, P. J., Smith, J. E., Heatherley, S. V., & Pleydell-Pearce, C. W. (2008). Time for tea: mood, blood pressure and cognitive performance effects of caffeine and theanine administered alone

and together. Psychopharmacology, 195(4), 569. Nehlig, A. (2010). Is caffeine a cognitive enhancer?. Journal of Alzheimer's Disease, 20(s1), S85-S94.

1.3 Research questions

The research question that will be examined in the experiment (or alternatively the hypothesis that will be tested in the experiment) Does coffee consumption increases IQ test score #Should we add the caffeine level ?

1.4 The related conceptual model

This model should include: *Independent variable(s)* -> *Coffee consumption* Dependent variable -> score at IQ test. *Mediating variable (at least 1)* -> *sleepingness feeling* Moderating variable (at least 1) -> amount of caffeine, prior coffee consumption habit/caffeine tolerance

1.5 Experimental Design

Note that the study should have a true experimental design The experiment is a two groups, post test only, randomized controlled trail.

1.6 Experimental procedure

Describe how the experiment will be executed step by step The participants will be separated into two groups randomly. One group will do the IQ test without any prior coffee consumption while the second group will do the test half an hour after coffee consumption. In the coffee consumption groups, participants will be separated in three subgroups where they will get coffee with different caffeine level. This will allow us to measure the general impact of drinking coffee on an IQ test but it will also allow us to test the difference between each caffeine level.

1.7 Measures

Describe the measure that will be used The Coffee consumption will be measured in ml. The performance in an IQ test will in a simple integer number on the scale from 0-200 where the mean is around 100. Sleepingness will be given by the participants on the scale from 0-10 where 10 means the highest level of sleepingness. The amount of caffeine will be measured in mg. Prior coffee drinking habits will be given by participants. They will be asked how much coffee they typically drink on a normal day.

1.8 Participants

Describe which participants will recruit in the study and how they will be recruited

Since we are just going to make this experiment on the effects of coffee consumption on students at TU Delft we need to find participants from that group of people. Emails will be sent out to the student body explaining the theory of the experiments and willing volunteers asked to fill in a form. We will try to contact an external company of some sort to get some credit or coupons that we can give to participants as a reward for helping out.

1.9 Suggested statistical analyses

Describe the statistical test you suggest to carry out on the collected data We will use a one way Analysis of Variance (ANOVA) test between group. Indeed, since the IQ test follows a gaussian distribution, we just want to compare the mean of each group.

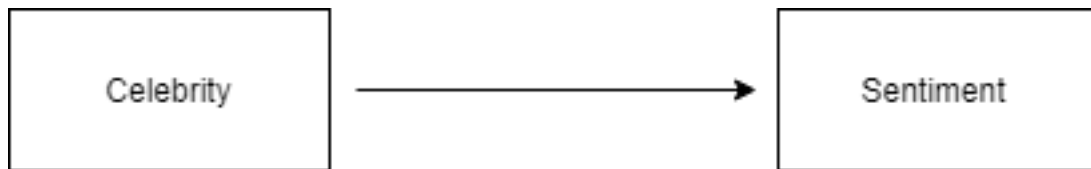


Figure 1: Conceptual model

2 Part 2 - Generalized linear models

2.1 Question 1 Twitter sentiment analysis (Between groups - single factor)

2.1.1 Collecting tweets, and data preparation

We collected Tweets for the three celebrities Beyonce, Madonna and Kanye West. The code can be found in the markdown file.

2.1.2 Conceptual model

Make a conceptual model for the following research question: Is there a difference in the sentiment of the tweets related to the different celebrities?

We can see that the sentiment of tweets related to different celebrity is directly connected to the celebrity itself. Therefore the conceptual model is very simple consisting of two variables, "Celebrity" and "Sentiment".

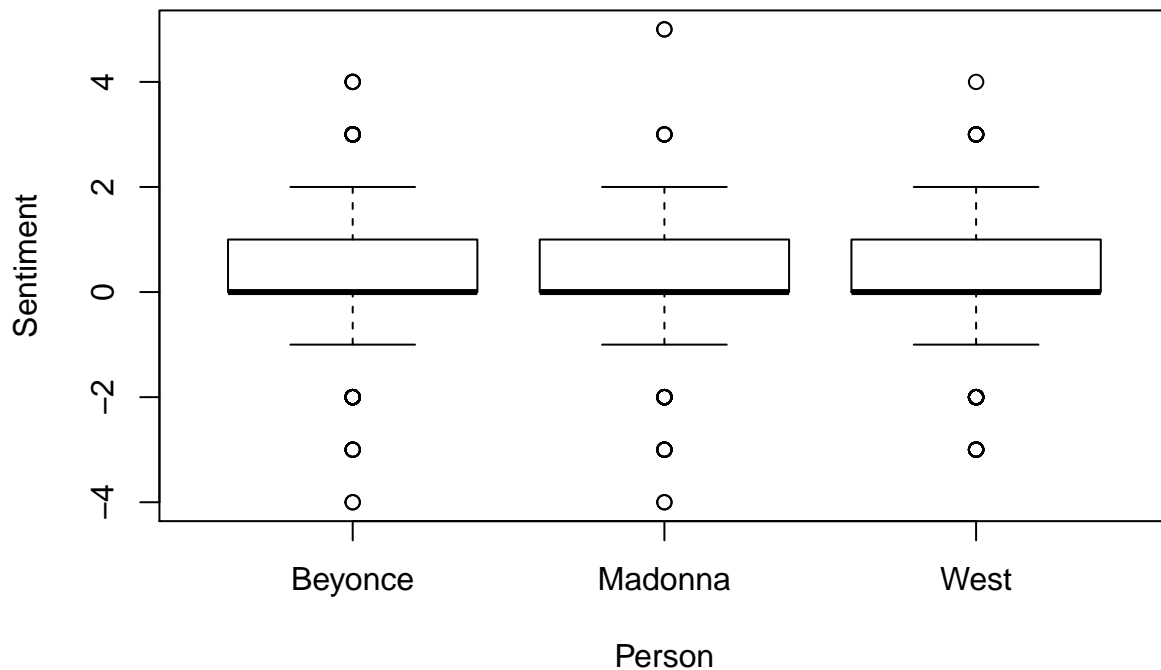
2.1.3 Homogeneity of variance analysis

Analyze the homogeneity of variance of sentiments of the tweets of the different celebrities

Lets start by looking at how the boxplot looks for each person and the relevant sentiment that has been analysed.

```
#this was not here in the intermediate report.  
#include your code and output in the document  
boxplot(score ~ Person, data=semFrame, main="Boxplot of sentiment for each person",  
        xlab="Person", ylab="Sentiment")
```

Boxplot of sentiment for each person



```
leveneTest( semFrame$score,semFrame$Person, center = median)
```

```
## Levene's Test for Homogeneity of Variance (center = median)
##      Df F value    Pr(>F)
## group  2  16.178 1.027e-07 ***
##      2997
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The Levene test results in a very low p-value., therefore the hypothesis of equal variances is rejected and it is concluded that there is a difference between the variances in the pipulation. Therefore the variance is not considered to be homogeneous.

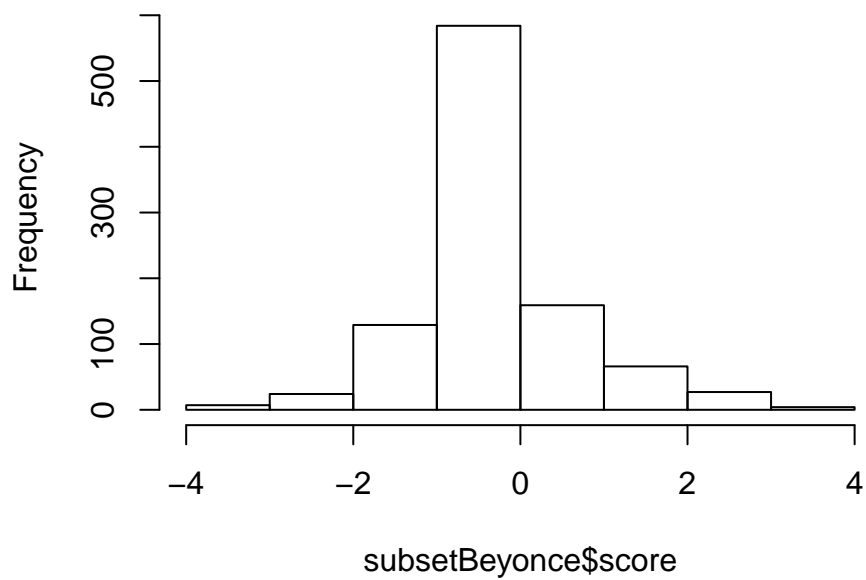
2.1.4 Visual inspection

Graphically examine the variation in tweets' sentiments for each celebrity (e.g. histogram, density plot)

#include your code and output in the document

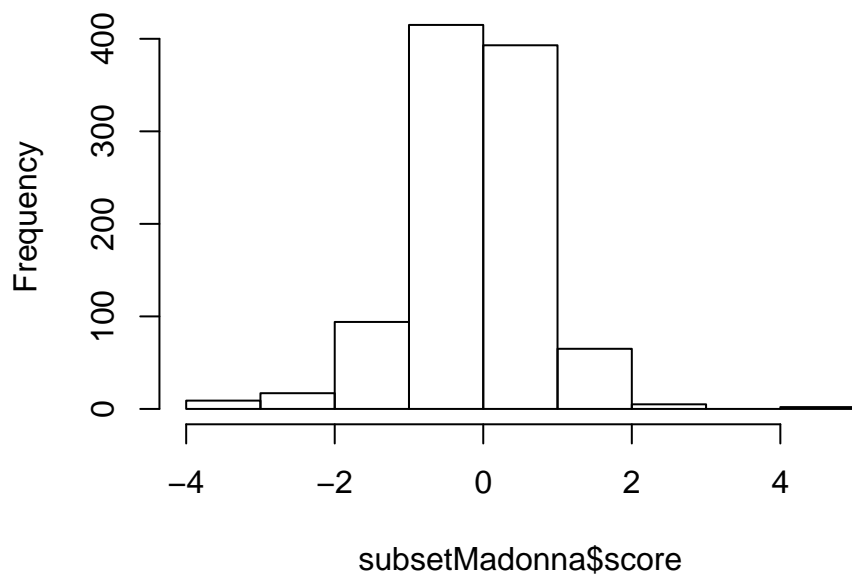
```
hist(subsetBeyonce$score)
```

Histogram of subsetBeyonce\$score



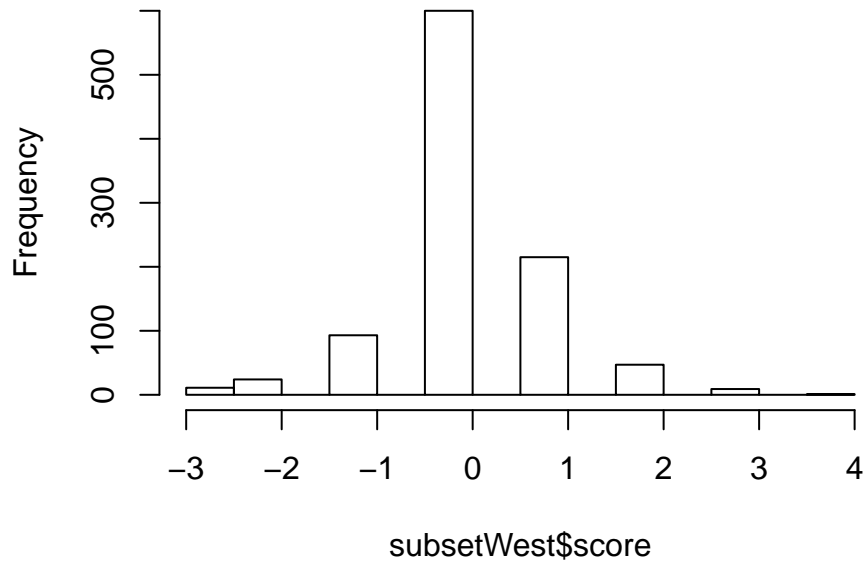
```
hist(subsetMadonna$score)
```

Histogram of subsetMadonna\$score



```
hist(subsetWest$score)
```

Histogram of subsetWest\$score

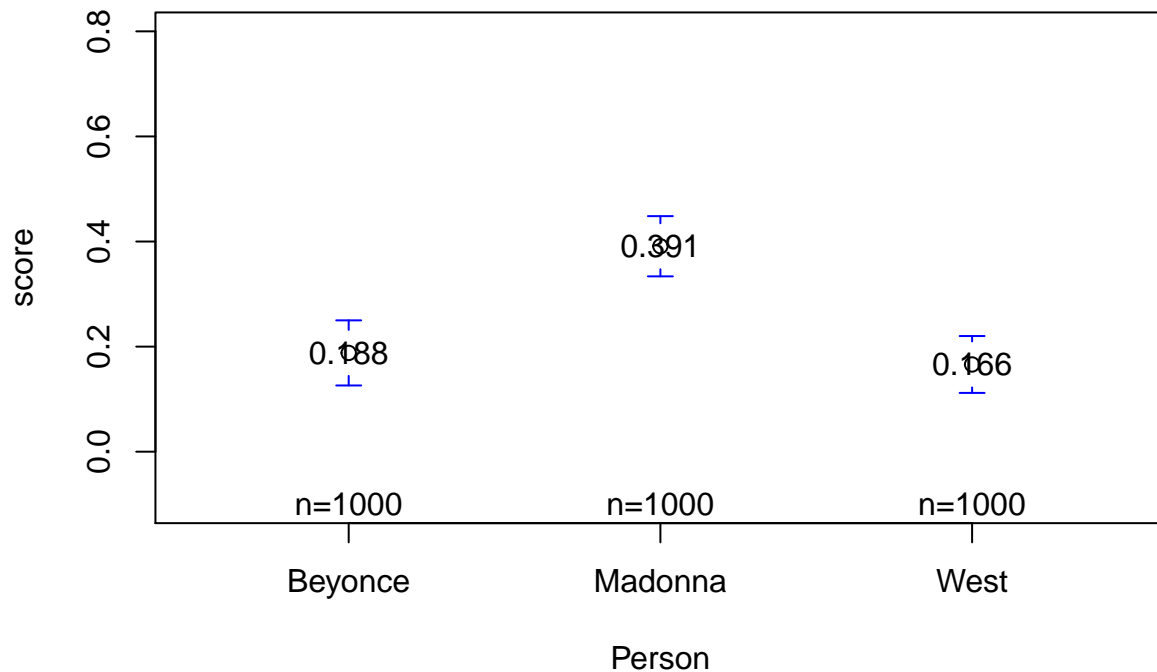


```
meanBeyonce<-mean(subsetBeyonce$score)
meanMadonna<-mean(subsetMadonna$score)
meanWest<-mean(subsetWest$score)
```

2.1.5 Mean sentiments

Here below we plot the means of each class using `plotmeans` from the package `gplots`. We can see that the mean for Beyonce is 0.188, for Madonna is 0.391 and for Kanye West it is 0.166. Where a lower value means that the sentiment analysis is more negative.

```
plotmeans(score ~ Person, data = semFrame, mean.labels = TRUE, connect = FALSE, ylim = c(-0.1, 0.8))
```



2.1.6 Linear model

```
#include your code and output in the document
model0<- lm(score ~ 1, data = semFrame) #model without predictor
model1<- lm(score ~ Person, data = semFrame) #model with predictor
AnovaResults <- anova(model0,model1)
```

The calculated f-value, $F(2,2997)$ is 17.6730435 and the p-value is 2.3420845×10^{-8} . Since the p-value is so small we can assume that the sentiment of tweets is significantly different depending on what celebrity is mentioned in the tweet.

2.1.7 Post Hoc analysis

If a model that includes the celebrity is better in explaining the sentiments of tweets than a model without such predictor, conduct a post-hoc analysis with e.g. Bonferroni correction, to examine which of celebrity tweets differ from the other celebrity tweets

```
#include your code and output in the document
BonferroniResults <- pairwise.t.test(semFrame$score, semFrame$Person, paired = FALSE, p.adjust.method = "bonferroni")
BonferroniP <- BonferroniResults$p.value
```

```
##           Beyonce      Madonna
## Madonna 3.613818e-06          NA
## West    1.000000e+00 2.251845e-07
```

We chose to use the Bonferroni correction to conduct this post-hoc analysis. There the p-values are multiplied by the number of comparisons.

According to results here above: Madonna vs West has significance that is adjusted to bonferroni of 2.2518449×10^{-7} , Madonna vs Beyonce has significance that is adjusted to bonferroni of 3.6138176×10^{-6} and Beyonce vs West has significance that is adjusted to bonferroni of 1. We can see that the bonferroni adjustment leads to all the p-values to be higher than the p-value obtained in the Anova test above. Since the value for Beyonce and Madonna is rather close to 1 it tells us that there really is not much difference on the sentiment for the two. But the other two pairs still have a significant difference, that is the sentiment is significantly different between the two persons in the pair.

2.1.8 Report section for a scientific publication

According to the analysis we have seen that the sentiment of Tweets is not always significantly correlated to the celebrity.

First we visualized the results a bit by looking at the plot of the means of the tweet sentiment per person. We saw that the means for Madonna and Beyonce are closest to each other.

A linear model was fitted on the number of the sentiment score, comparing the difference when taking the relative celebrity in account and not. We first conducted an Anova test and obtained the p value of 2.3420845×10^{-8} which hinted there was significant difference in the sentiment, depending on which celebrity it is for. Then A Post Hoc analysis by the means of Bonferroni was conducted. There we saw that the difference of sentiment is not significant for Beyonce and Madonna (p-value of 3.6138176×10^{-6}), but is significant for Kanye West vs Madonna and Beyonce (p-values: $1, 2.2518449 \times 10^{-7}$).