# Report coursework assignment A - 2021
## CS4125 Seminar Research Methodology for Data Science

Nikki Bouman (4597648), Anuj Singh (), Gwennan Smitskamp ()

20/04/2021

## Contents

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

## 1.1 The motivation for the planned research

(Max 250 words)

## 1.2 The theory underlying the research

(Max 250 words) Preferable based on theories reported in literature

## 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)

## 1.4 The related conceptual model

This model should include: *Independent variable(s)* Dependent variable *Mediating variable (at least 1)* Moderating variable (at least 1)

## 1.5 Experimental Design

Note that the study should have a true experimental design

## 1.6 Experimental procedure

Describe how the experiment will be executed step by step

## 1.7 Measures

Describe the measure that will be used

## 1.8 Participants

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

## 1.9 Suggested statistical analyses

Describe the statistical test you suggest to care out on the collected data

# 2 Part 2 - Generalized linear models

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

### 2.1.1 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?

### 2.1.2 Collecting tweets, and data preparation

Include the annotated R script (excluding your personal Keys and Access Tokens information), but put echo=FALSE, so code is not included in the output pdf file.

### 2.1.3 Homogeneity of variance analysis

Analyze the homogeneity of variance of sentiments of the tweets of the different celebrities, and provide interpretation
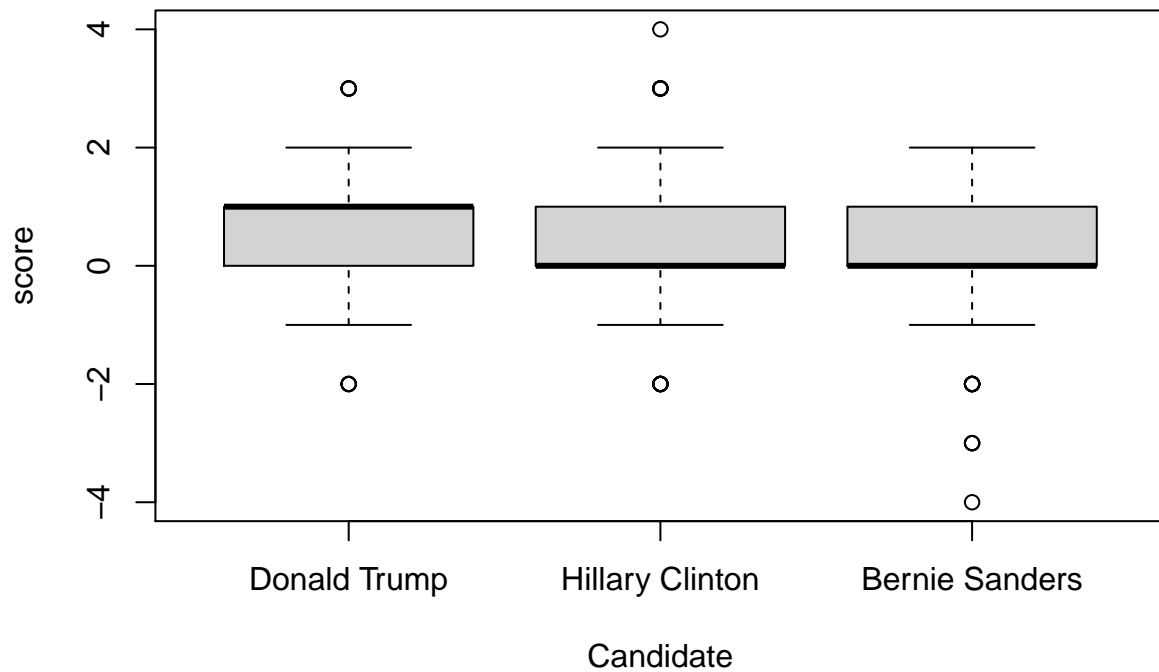
```
#include your code and output in the document
leveneTest(semFrame$score, semFrame$Candidate, center = median)

## Levene's Test for Homogeneity of Variance (center = median)
##        Df F value Pr(>F)
## group   2  0.3666 0.6933
##       393
```

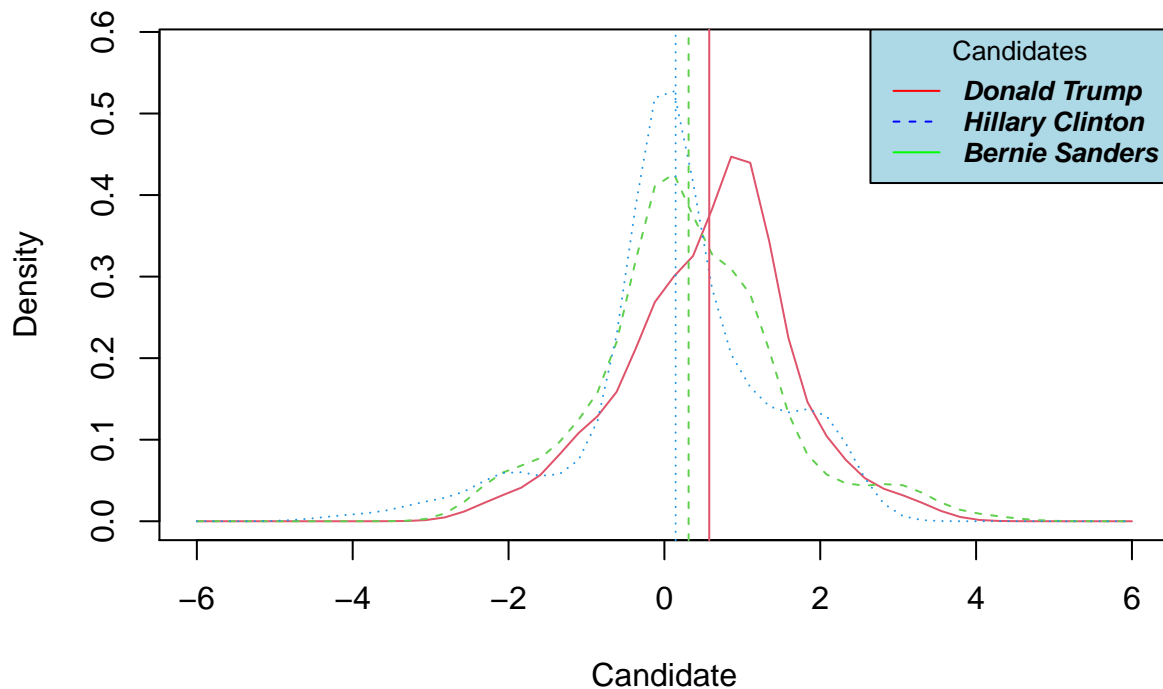### 2.1.4 Visual inspection Mean and distribution sentiments

Graphically examine the mean and distribution sentiments of tweets for each celebrity, and provide interpretation

```
boxplot(score ~ Candidate, data = semFrame)
```

```
sm.density.compare(semFrame$score, semFrame$Candidate, xlab = "Candidate")
  title(main="Visual inspection Mean and distribution sentiments")
  legend('topright', legend = levels(semFrame$Candidate), col=c('red', 'blue', 'green'),
         title="Candidates", lty=1:2, cex=0.8, text.font = 4, bg='lightblue')
means <- aggregate(semFrame$score ~ semFrame$Candidate, FUN = mean)
abline(v = means[1,2], col = 2)
abline(v = means[2,2], col = 3, lty = 2)
abline(v = means[3,2], col = 4, lty = 3)
```

## Visual inspection Mean and distribution sentiments

### 2.1.5 Frequentist approach

**2.1.5.1 Linear model** Use a linear model to analyze whether the knowledge to which celebrity a tweet relates has a significant impact on explaining the sentiments of the tweets. Provide interpretation of results

```
model0 <- lm(formula = score ~ 1 , data = semFrame)
model1 <- lm(formula = score ~ Candidate , data = semFrame)
anova(model0, model1, test = "F")
```

```
## Analysis of Variance Table
##
## Model 1: score ~ 1
## Model 2: score ~ Candidate
##   Res.Df    RSS Df Sum of Sq      F   Pr(>F)
## 1    395 481.29
## 2    393 468.77  2     12.52 5.2482 0.005632 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
summary(model1)
```

```
##
## Call:
## lm(formula = score ~ Candidate, data = semFrame)
##
## Residuals:
```

```
##     Min      1Q  Median      3Q     Max
## -4.1439 -0.3106 -0.1439  0.6894  3.6894
##
## Coefficients:
##                         Estimate Std. Error t value Pr(>|t|)
## (Intercept)              0.57576    0.09506   6.057 3.25e-09 ***
## CandidateHillary Clinton -0.26515    0.13444  -1.972  0.04927 *
## CandidateBernie Sanders  -0.43182    0.13444  -3.212  0.00143 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.092 on 393 degrees of freedom
## Multiple R-squared:  0.02601,    Adjusted R-squared:  0.02106
## F-statistic: 5.248 on 2 and 393 DF,  p-value: 0.005632
```

```
smodel0 <-summary(model0)
llm0 <-sum(dnorm(semFrame$score, mean = predict(model0), sd= smodel0$sigma, log=TRUE))
AIC_model0 <- -2*llm0 + 2*2
AIC_model0
```

```
## [1] 1205.045
```

```
smodel1 <-summary(model1)
llm1 <-sum(dnorm(semFrame$score, mean = predict(model1), sd= smodel1$sigma, log=TRUE))
AIC_model1 <- -2*llm1 + 2*3
AIC_model1
```

```
## [1] 1196.617
```
```
#include your code and output in the document
```

**2.1.5.2  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. Provide interpretation of the results

```
pairwise.t.test(semFrame$score, semFrame$Candidate, paired = FALSE, p.adjust.method = "bonferroni")
```

```
##
##  Pairwise comparisons using t tests with pooled SD
##
## data:  semFrame$score and semFrame$Candidate
##
##                 Donald Trump Hillary Clinton
## Hillary Clinton 0.1478       -
## Bernie Sanders  0.0043       0.6474
##
## P value adjustment method: bonferroni
```
```
#include your code and output in the document
```

**2.1.5.3  Report section for a scientific publication**  Write a small section for a scientific publication, in which you report the results of the analyses, and explain the conclusions that can be drawn.

**2.1.6  Bayesian Approach**

**2.1.6.1  Model description**  Describe the mathematical model fitted on the most extensive model. (hint, look at the mark down file of the lectures to see example on formulate mathematical models in markdown).

Justify the priors.

Comparing multiple levels

**2.1.6.2   Model comparison**   Conduct model analysis and provide brief interpretation of the results

```
#semFrame1 <-subset(semFrame, (Candidate == "Donald Trump"))
#semFrame2 <-subset(semFrame, (Candidate == "Hillary Clinton"))
#semFrame3 <-subset(semFrame, (Candidate == "Bernie Sanders"))
#fit <-bayes.t.test(semFrame1$score,semFrame2$score)
#show(fit)
#plot(fit)

#fit <- bayes.t.test(semFrame$score)
#plot(fit)

#semFrame$yearF <-factor(semFrame$year, levels = c(2004:2007), labels = c(2004:2007))
#da <-subset(semFrame, select = c(t_term1, yearF))

m0 <-map2stan(alist(score ~ dnorm(mu, sigma),mu <-a ,a ~ dnorm(50, 25),sigma ~ dunif(0.001, 40)),  data
```

```
## Computing WAIC
```

```
m1 <-map2stan(alist(score ~ dnorm(mu, sigma),mu <-a[Candidate] ,a[Candidate] ~ dnorm(50, 25),sigma ~ du
```

```
## Computing WAIC
```

```
compare(m0, m1, func=WAIC)
```

```
##          WAIC      SE    dWAIC      dSE    pWAIC     weight
## m1 1199.145 35.21491 0.000000       NA 4.479862 0.96408731
## m0 1205.725 35.46006 6.580182 6.199231 2.610444 0.03591269
```

```
precis(m1, depth = 2, prob = .95)
```

```
##             mean         sd        2.5%      97.5%     n_eff      Rhat4
## a[1]   0.5762010 0.09460952  0.38910477  0.7606644 24711.99 1.0000105
## a[2]   0.3102428 0.09530952  0.12505815  0.4963903 24420.81 1.0000818
## a[3]   0.1443300 0.09481419 -0.04398948  0.3295576 23594.74 0.9999367
## sigma 1.0958408 0.03924441  1.02286058  1.1758316 24578.47 0.9999503
```

```
#include your code and output in the document
```

**2.1.6.3   Comparison celebrity pair**   Compare sentiments of celebrity pairs and provide a brief interpretation (e.g. CIs)