

Report coursework assignment A - 2021

CS4125 Seminar Research Methodology for Data Science

Nikki Bouman (4597648), Anuj Singh (5305926), Gwennan Smitskamp (4349822)

12/05/2021

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	3
1.7	Measures	3
1.8	Participants	3
1.9	Suggested statistical analyses	5
2	Part 2 - Generalized linear models	5
2.1	Question 1 Twitter sentiment analysis (Between groups - single factor)	5
2.1.1	Conceptual model	5
2.1.2	Collecting tweets, and data preparation	5
2.1.3	Homogeneity of variance analysis	6
2.1.4	Visual inspection Mean and distribution sentiments	6
2.1.5	Frequentist approach	7
2.1.6	Bayesian Approach	9
3	Part 3 - Multilevel model	10
3.1	Visual inspection	10
3.2	Frequentist approach	13
3.2.1	Multilevel analysis	13
3.2.2	Report section for a scientific publication	16
3.3	Bayesian approach	17
3.3.1	Model description	17
3.3.2	Model comparison	17
3.3.3	Estimates examination	18

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

1.1 The motivation for the planned research

The recent outbreak of the COVID-19 pandemic has changed our daily lives significantly. People are obligated to stay at home, disrupting their usual social interactions in both work and private life. The situation compels people to meet online. Typically, in such digital interactions, interlocutors can see each other by means of

webcam streaming. However, this may not always be the case. Some or all interlocutors may not be visible during online dialogue, which could affect the quality of the conversation and the mutual understanding.

An important effect of the shift from face-to-face to online interaction can be revealed by studying laughter, as it is extremely contagious social behavior (Provine, 1992). Humans are very prone to unintentionally or unconsciously laugh as a social signal in any form; from a minor smile to laughing out loud. Additionally, laughing is one of the most important social signals for lubricating the flow of social interaction (Griffin et al., 2015).

1.2 The theory underlying the research

The effect of visibility on the use of gestures as a communicative function has been studied broadly (Alibali, Heath, & Myers, 2001; J. B. Bavelas, Chovil, Lawrie, & Wade, 1992; Cohen & Harrison, 1973; Cohen, 1977; Emmorey & Casey, 2001; Krauss, Dushay, Chen, & Rauscher, 1995; Rimé, 1982). J. Bavelas, Gerwing, Sutton, and Prevost (2008) provide a summary of previous experiments where rate and form of gestures were compared under two conditions: where the addressee could see the speaker and where the addressee could not see the speaker. These experiments show that speakers gestured at higher rate when they communicated with mutual visibility than without. J. Bavelas et al. (2008) extended these experiments by focusing on both visibility and dialogue as a variable, finding similar results. Furthermore, they found that speakers gestured at a significantly higher rate in a telephone dialogue than in a monologue to a tape recorder, confirming that visibility is not the only variable operating in telephone conversations. These experiments showed us that visibility plays a major role in the rate of gesturing, but that people also gesture when they are not visible to each other. As laughter can be seen as a form of gesturing, these findings are relevant for this study.

Laughing together is found to be essentially collaborative (Mehu & Dunbar, 2008; Coates, 2007). Joint laughter therefore serves important means to achieve effective team meetings (Ponton, Osbourne, Greenwood, & Thompson, 2018), considering that people who laugh on video are perceived with a higher likeability than people who do not (Reysen, 2006). This social function of joint laughter emphasises the relevance of studying the occurrence, now that the majority of meetings take place online.

1.3 Research questions

We will aim to answer the following research question:

What is the effect of webcam visibility during online dialogue on the frequency of joint laughter?

When recognizing laughter we do not focus on the reason why someone is laughing. We consider anything from an awkward laugh in a moment of silence to laughing out loud about a joke as a laughter episode regardless of the context.

1.4 The related conceptual model

The conceptual model related to the research question can be viewed in Figure 1. The main question is about the effect of mutual visibility on joint laughter. The mediating variable is familiarity, which we can define as the level of friendliness or intimacy between people. This can be caused when people know each other (which we aim to avoid), but also when people find similarities in their interests and behaviours. It is even possible that people of the same gender will feel more familiar with each other. A moderating variable is the duration of the experiment. The participants will do the experiment for about one hour, which can have a negative effect on the frequency of laughter, thus on the frequency of joint laughter.

1.5 Experimental Design

One of the most important requirements for the setup of the experiment, was the creation of a comfortable and pleasant ambiance so that people would laugh. Therefore, it was decided that a game would comply, as the participants get the chance to interact with each other in a undemanding setting where the attention of the participants would be drawn to a task. It was reasoned that this would contribute to a reduction of

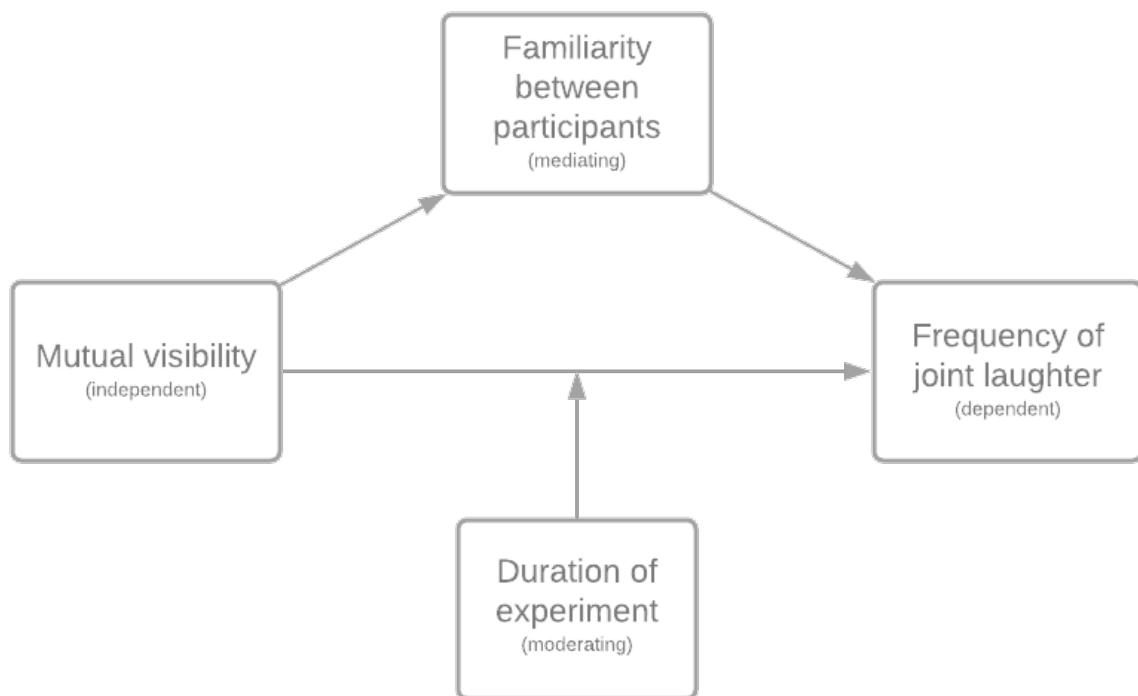


Figure 1: The conceptual model to test the effect of mutual webcam visibility on joint laughter

awkwardness and give all the participants the option to speak and laugh. Additionally, the game needed to have a smooth flow that would automatically keep going to keep the interference of the researchers to a minimum.

The game that was chosen is called 30 Seconds. During the game, participants work together in teams (in the case of the experiment: two teams of two people) and gain points by guessing what the team member is describing. These descriptions include concepts such as famous persons, locations, movies and brands. Every team gets 30 seconds to guess as many concepts on the card as possible. Who is describing and who is guessing switches after every card.

1.6 Experimental procedure

The experiment will be set up in an online setting in a Zoom meeting. The host, one of us (not visible), will be able to send private messages containing the five words, share their screen and sound for a 30 second timer, and turn the participants' webcams on and off.

The following will repeat for every card of five words:

1. The host sends the words to the participant who has the turn to describe.
2. The host starts the 30 second timer.
3. The participant will try to describe as many words as possible, while his/her teammate will try to guess the words.
4. The timer rings, the host puts the score in the chat.

During the experiment, multiple things will happen. Each time all players have guessed and described a card (i.e. after four cards total), their webcams will switch on or off. After each player has guessed and described four cards (i.e. after sixteen cards total), the final score will be displayed and the teams will be rearranged. The previous will then repeat until every participant has been in a team with every other participant. An example of such an experiment is displayed in Figure 2. To counterbalance the experiments, half of the experiments will start with the webcams on, while the other half will start with the webcams off.

1.7 Measures

The data that has been collected includes audio and video of participants. The first step in data analysis involves annotating the signals. This will be done with a program in which we can manually select timesteps in which the participant is laughing. In the end we will have annotations for every person that contains the total amount of laughter (frequency), the amount of laughter with their webcam off, and the amount of laughter with their webcam on..

In a case where different people will annotate this data, we will first let every one of them annotate the same data sample. Then we can calculate the consistency in the annotations with for example Krippendorff's Alpha or Cohen's Kappa. When this is high enough, they can start annotating separate data.

1.8 Participants

The research is not specifically about a certain group of people, but more in general. However, we do want to aim for people with experience in an online setting and people who speak the same language. To be exact, we will conduct the experiment with people from the ages of 18 to 50 who speak dutch. The number of participants depends on the acceptable margin or error, which we do not know since we would have to go much more in-depth. However, since the population size basically includes more than 10.000 people and the independent variable is categorical, we should aim for around 385 participants (W.P. Brinkman, 2009). This results into around 96 experiments.

The experiments should be easy to conduct since the participants participate online; there is no need to travel. Moreover, the whole experiment should not take long. There will be 48 cards and for each card they

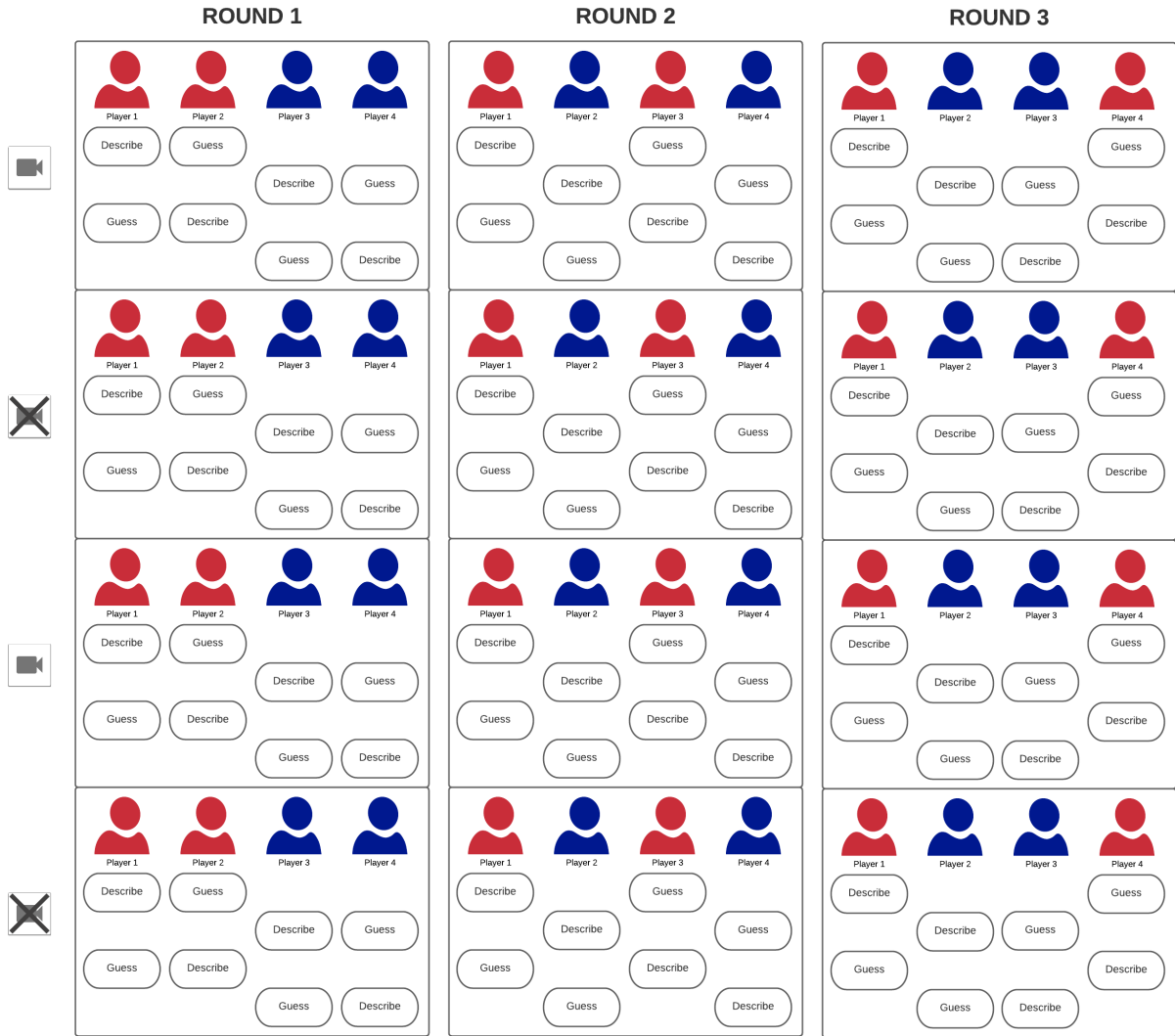


Figure 2: An example of an experimental setup.

have 30 seconds to explain. Taking some talking afterwards into account, the experiment should take 45 to 60 minutes.

To find the participants we can thus search online. Using a medium on which we state the experiment, we can find dutch people from all around the Netherlands of different age groups. They could for example sign up for a specific date, and if four people have signed up, the experiment can be held.

1.9 Suggested statistical analyses

To determine the significance of the results we will subject the data to statistical tests. To test the effect of mutual visibility (categorical, since it is either on or off) on the frequency (numerical) we will use the Wilcoxon signed rank test. The Wilcoxon signed rank test is a suitable statistical test when the measurements for a single variable are taken under two different conditions. It is similar to the paired t-test, however the paired t-test assumes that the data is normally distributed which we cannot assume for frequency of laughter. It tests the null hypothesis that the median difference of a two sets of observations is zero.

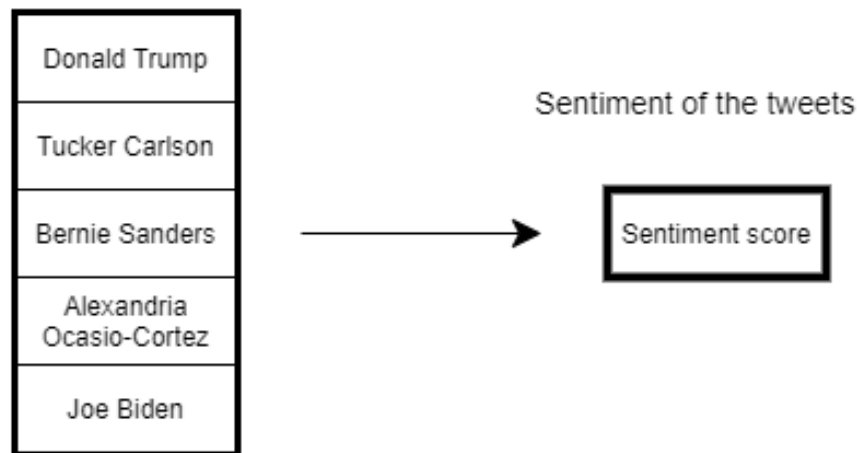
2 Part 2 - Generalized linear models

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?

Relation to different celebrities



Is there a difference in the sentiment of the tweet related to the different celebrities?

Figure 3: The conceptual model for the 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

We found five celebrities in US politics: Donald Trump, Tucker Carlson, Bernie Sanders, Alexandria Ocasio-Cortez, Joe Biden. As dutch students we are not well-versed in the popular English twitter celebrities, so US politics was the best option for us to find celebrities that had enough recent Tweets for the Twitter API.

2.1.3 Homogeneity of variance analysis

```
pander(leveneTest(semFrame$score, semFrame$Celeb))
```

Table 1: Levene's Test for Homogeneity of Variance (center = median)

	Df	F value	Pr(>F)
group	4	23.64	5.714e-19
	1465	NA	NA

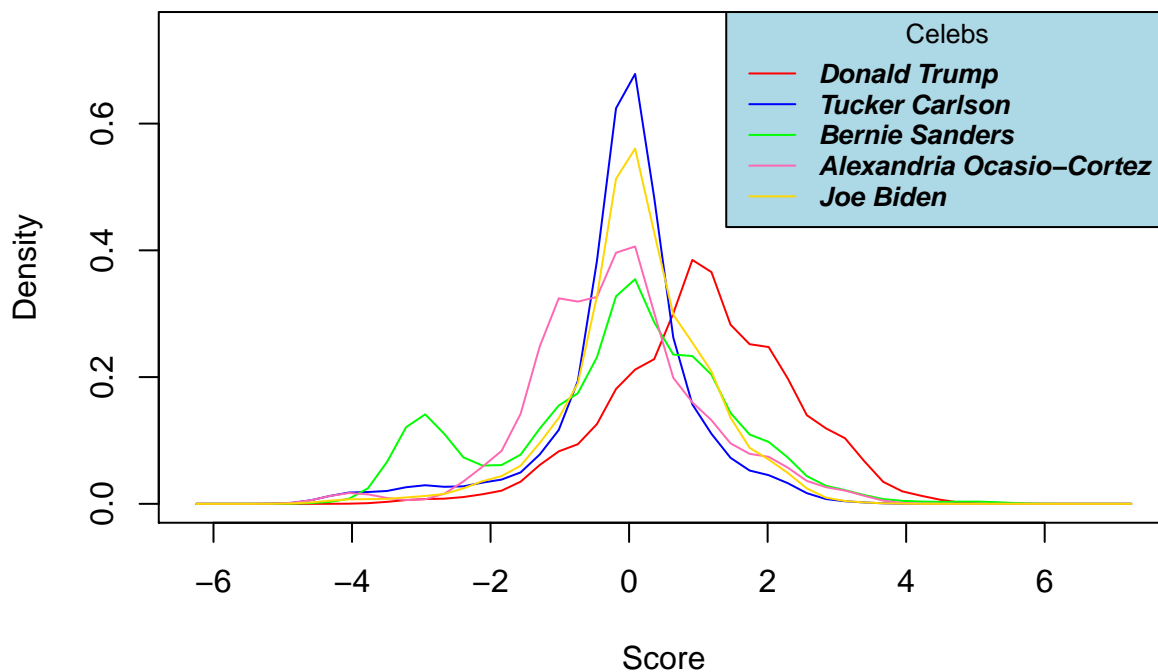
The Levene test reveals a p-value smaller than 0.001, indicating that there is significant difference between the group variances in sentiment score. We conclude that the variance among the five groups is not equal.

2.1.4 Visual inspection Mean and distribution sentiments

We plot both a line density and a distribution histogram which includes a mean line.

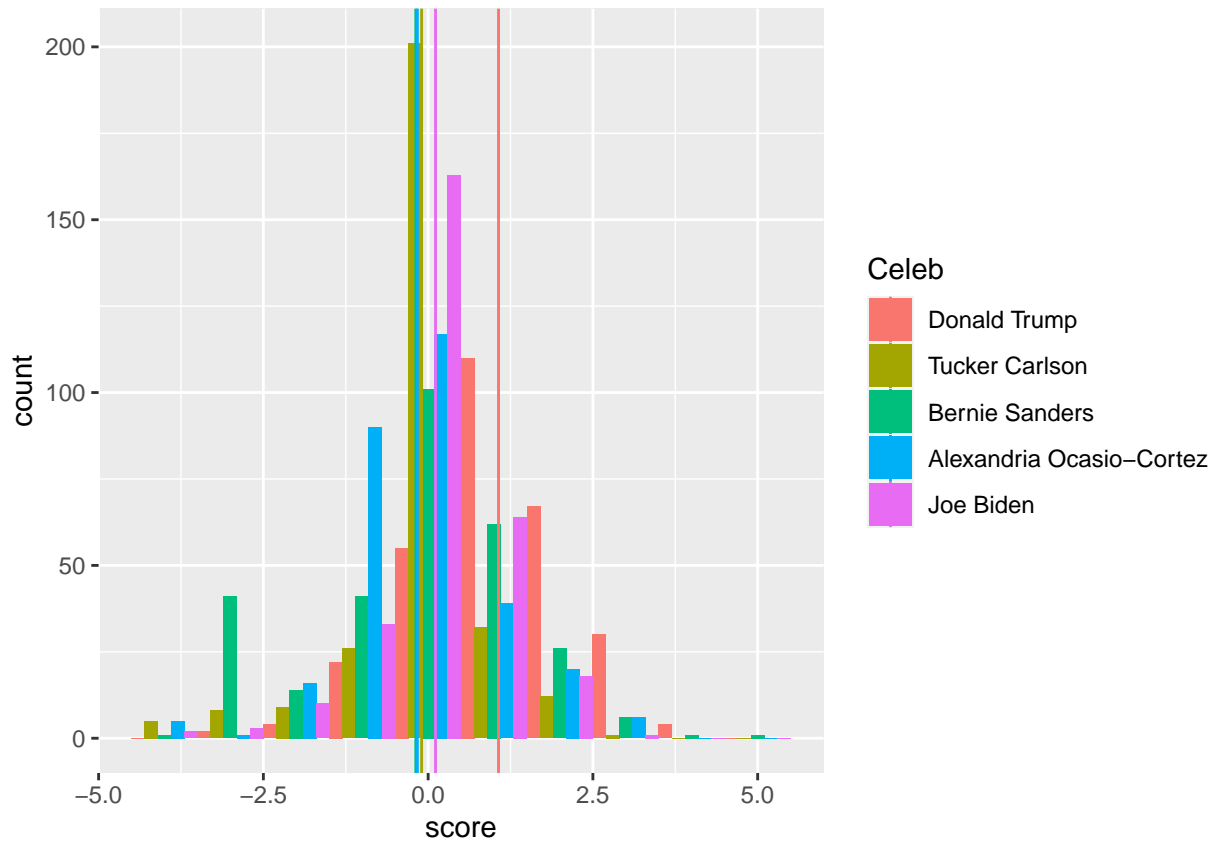
```
#boxplot(score ~ Celeb, data = semFrame, text.font = 4)
sm.density.compare(semFrame$score, semFrame$Celeb, xlab = "Score",
                   col=c('red', 'blue', 'green', 'hotpink', 'gold'), lty=c(1,1,1,1,1))
title(main="Visual inspection Mean and distribution sentiments")
legend('topright', legend = levels(semFrame$Celeb),
       col=c('red', 'blue', 'green', 'hotpink', 'gold'),
       title="Celebs", lty=1, cex=0.8, text.font = 4, bg='lightblue')
```

Visual inspection Mean and distribution sentiments



```
cdat <- ddply(semFrame, "Celeb", summarise, score.mean=mean(score))
ggplot(semFrame, aes(x=score, fill=Celeb)) +
  geom_histogram(binwidth=1, position="dodge") +
```

```
geom_vline(data=cdat, aes(xintercept=score.mean, colour=Celeb),
           linetype="solid", size=0.5)
```



We see all US politic Celebs have a mean around 0. #trump has the highest sentiment mean and the largest difference with the rest.

2.1.5 Frequentist approach

2.1.5.1 Linear model A one-way between subjects ANOVA was conducted to compare the effect of relation to celebrities on sentiment score in five conditions.

```
modelTwitter0 <- lm(formula = score ~ 1, data = semFrame)
modelTwitter1 <- lm(formula = score ~ Celeb, data = semFrame)
pander(anova(modelTwitter0, modelTwitter1, test = "F"))
```

Table 2: Analysis of Variance Table

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1469	2476	NA	NA	NA	NA
1465	2147	4	328.6	56.06	4.859e-44

There was a significant effect of relation to celebrities on sentiment score at the $p < .001$ level for the five conditions [$F(4, 1465) = 56.06$, $p < 0.001$].

```
#AIC
modelsTwitter <- list(modelTwitter0, modelTwitter1)
modelTwitter.names <- c("modelTwitter0", "modelTwitter1")
```



```
pander(aictab(cand.set = modelTwitter, modnames=modelTwitter.names),
        caption="Model selection based on AICc.")
```

Table 3: Model selection based on AICc. (continued below)

	Modnames	K	AICc	Delta_AICc	ModelLik	AICcWt	LL
2	modelTwitter1	6	4741	0	1	1	-2364
1	modelTwitter0	2	4942	201.3	1.939e-44	1.939e-44	-2469

Cum.Wt	
2	1
1	1

A lower AIC indicates a better fit, which is the model with the predictor.

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

2.1.5.2 Post Hoc analysis

```
## Warning in pander.default(pairwise.t.test(semFrame$score, semFrame$Celeb, : No
## pander.method for "pairwise.htest", reverting to default.
```

- **method:** t tests with pooled SD
- **data.name:** *semFrame\$score* and *semFrame\$Celeb*
- **p.value:**

Table 5: Table continues below

	Donald Trump	Tucker Carlson	Bernie Sanders
Tucker Carlson	3.197e-29	NA	NA
Bernie Sanders	1.727e-33	1	NA
Alexandria Ocasio-Cortez	2.575e-32	1	1
Joe Biden	4.257e-20	0.3485	0.02765

Alexandria Ocasio-Cortez	
Tucker Carlson	NA
Bernie Sanders	NA
Alexandria Ocasio-Cortez	NA
Joe Biden	0.05864

- **p.adjust.method:** bonferroni

Post hoc comparisons using the Bonferroni correction indicated that the corrected p-value for the trump condition was significantly different than the other conditions ($p < 0.001$). However, between the others condition it does not show a significant difference.

2.1.5.3 Report section for a scientific publication A one-way between subjects ANOVA was conducted to compare the effect of relation to celebrities on sentiment score in five conditions. There was a significant effect of relation to celebrities on sentiment score at the $p < .001$ level for the five conditions [$F(4, 1465) = 56.06, p < 0.001$]. However, post hoc comparisons using the Bonferroni correction indicated that only the corrected p-value for the trump condition was significantly different than the other conditions ($p < 0.001$), between the others condition it does not show a significant difference. Taken together, these results suggest that some celebrities really do have an effect on the sentiment in Tweets.

2.1.6 Bayesian Approach

2.1.6.1 Model description The sentiment scores seem to center around 0, and all seem to be single digits.

$$\begin{aligned} score &\sim \text{Norm}(\mu, \sigma) \\ \mu &= \alpha + b * \text{Celeb} \\ \alpha &= \text{Norm}(0, 10) \\ \sigma &= \text{Uniform}(0.001, 10) \end{aligned}$$

```
mTwitter0 <-map2stan(alist(
  score ~ dnorm(mu, sigma),
  mu <-a,
  a ~ dnorm(0, 10),
  sigma ~ dunif(0.001, 10)),
  data = semFrame , iter= 10000, chains = 4, cores = 4 )
```

2.1.6.2 Model comparison

Computing WAIC

```
mTwitter1 <-map2stan(alist(
  score ~ dnorm(mu, sigma),
  mu <-a[Celeb] ,
  a[Celeb] ~ dnorm(0, 10),
  sigma ~ dunif(0.001, 10)),
  data = semFrame , iter= 10000, chains = 4, cores = 4 )
```

Computing WAIC

```
pander(compare(mTwitter0, mTwitter1, func=WAIC))
```

	WAIC	SE	dWAIC	dSE	pWAIC	weight
mTwitter1	4741	70.16	0	NA	6.674	1
mTwitter0	4942	67.72	201.2	28.91	2.543	2.072e-44

Lower WAIC indicates a better performing model, so with predictors (mTwitter1) is the winning model.

```
pander(precis(mTwitter1, depth=2, prob = .95))
```

2.1.6.3 Comparison celebrity pair

	mean	sd	2.5%	97.5%	n_eff	Rhat4
a[1]	1.068	0.06997	0.9305	1.206	31991	0.9999
a[2]	-0.09853	0.07112	-0.2375	0.04129	32702	0.9999
a[3]	-0.1868	0.07076	-0.3268	-0.04858	29425	1
a[4]	-0.1635	0.07114	-0.3028	-0.02381	31660	0.9999
a[5]	0.112	0.07037	-0.02546	0.2503	31797	0.9999
sigma	1.212	0.02253	1.169	1.257	31995	1

Looking at the credibility intervals of the celebrities effects, We see the conditions where the mean of a condition does not fall within a credibility interval of an other condition. This holds for the a[1] (Trump) condition and a couple other combinations. We can again conclude that some celebrities really do have an effect on the sentiment in Tweets.

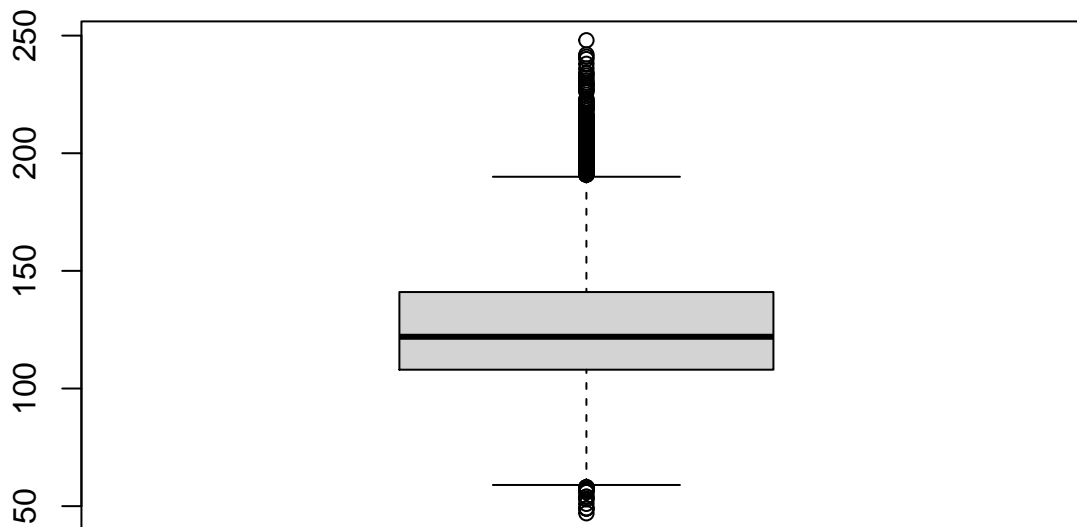
3 Part 3 - Multilevel model

3.1 Visual inspection

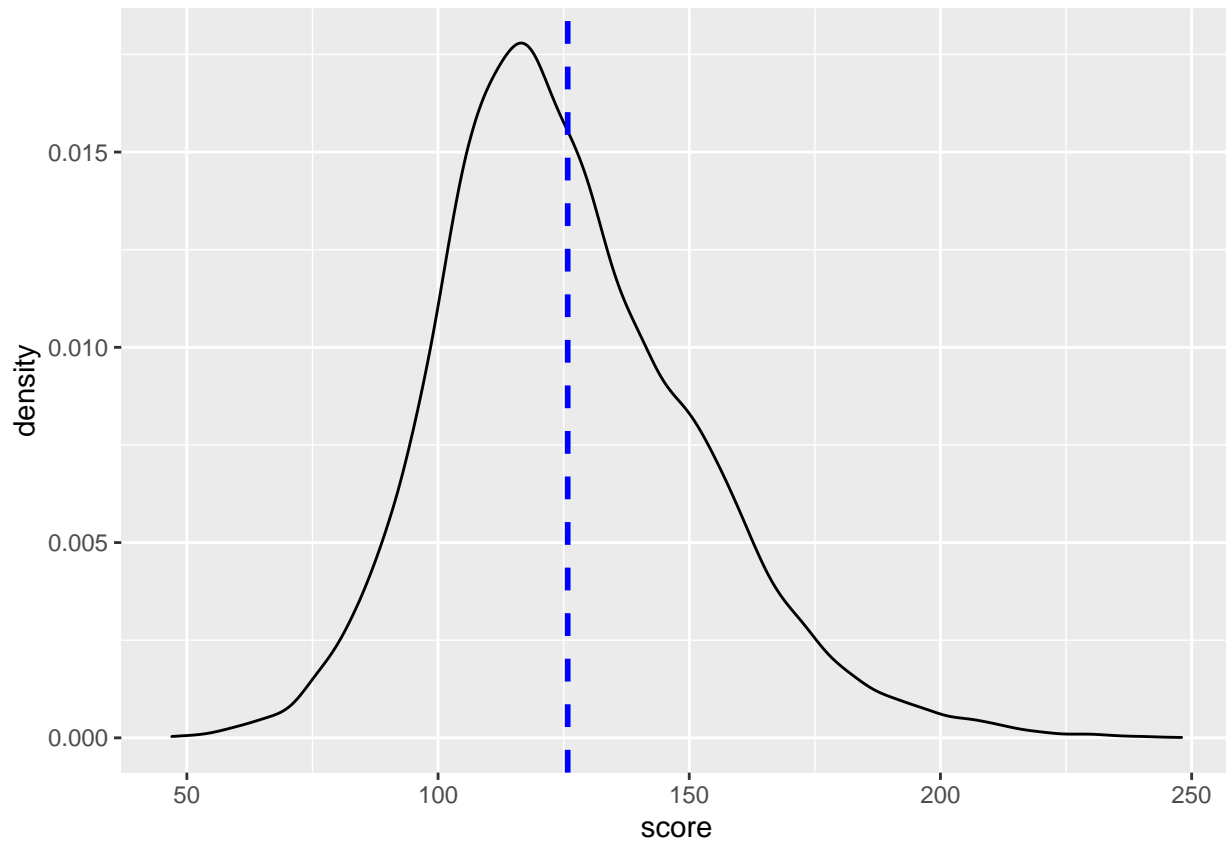
The boxplot and density plot show the distribution of the score. We can see that the mean score is 122 points. The minimum is set at 59, with outliers until 46, while the maximum is set at 190, with outliers until 248.

```
# Get data
filepath <- ("set0.csv")
ds <- read.csv(file=filepath, header=TRUE)
ds <- data.frame(ds)

# boxplot score overall distribution (session independent)
boxplot(ds$score)
```



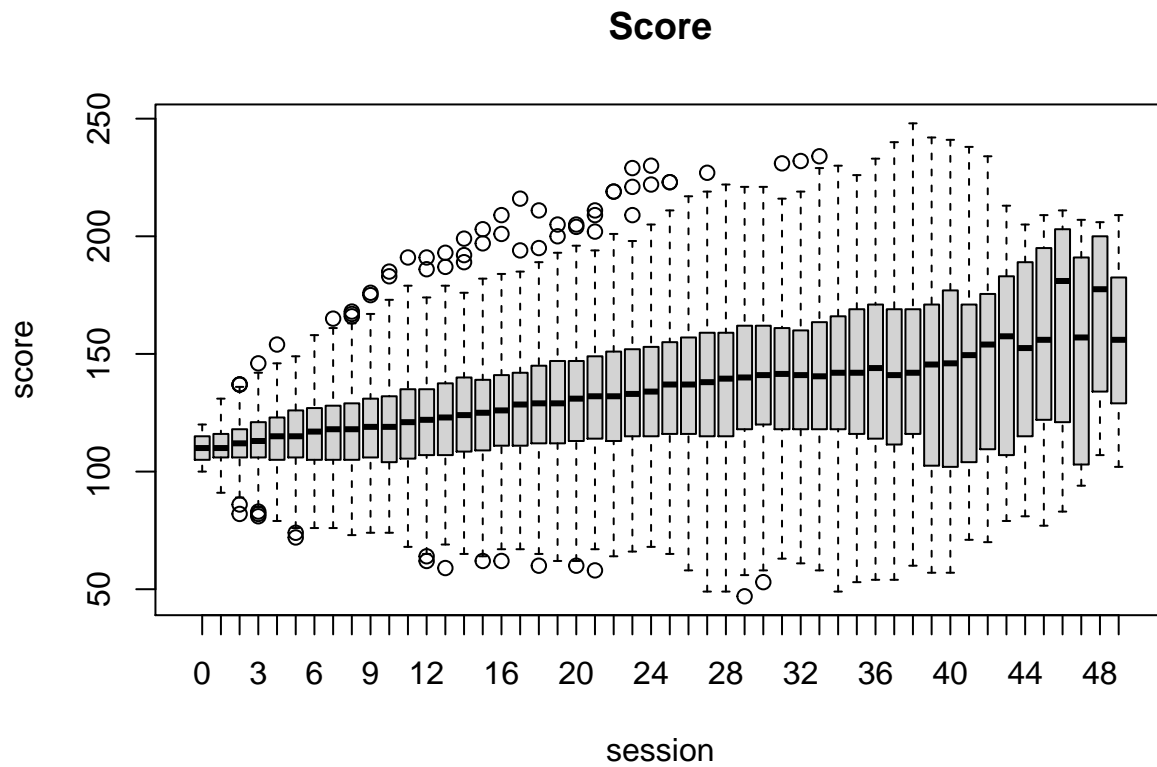
```
# density score overall distribution (with mean line)
p <- ggplot(ds, aes(x=score)) + geom_density()
p + geom_vline(aes(xintercept=mean(score)), color="blue", linetype="dashed", size=1)
```



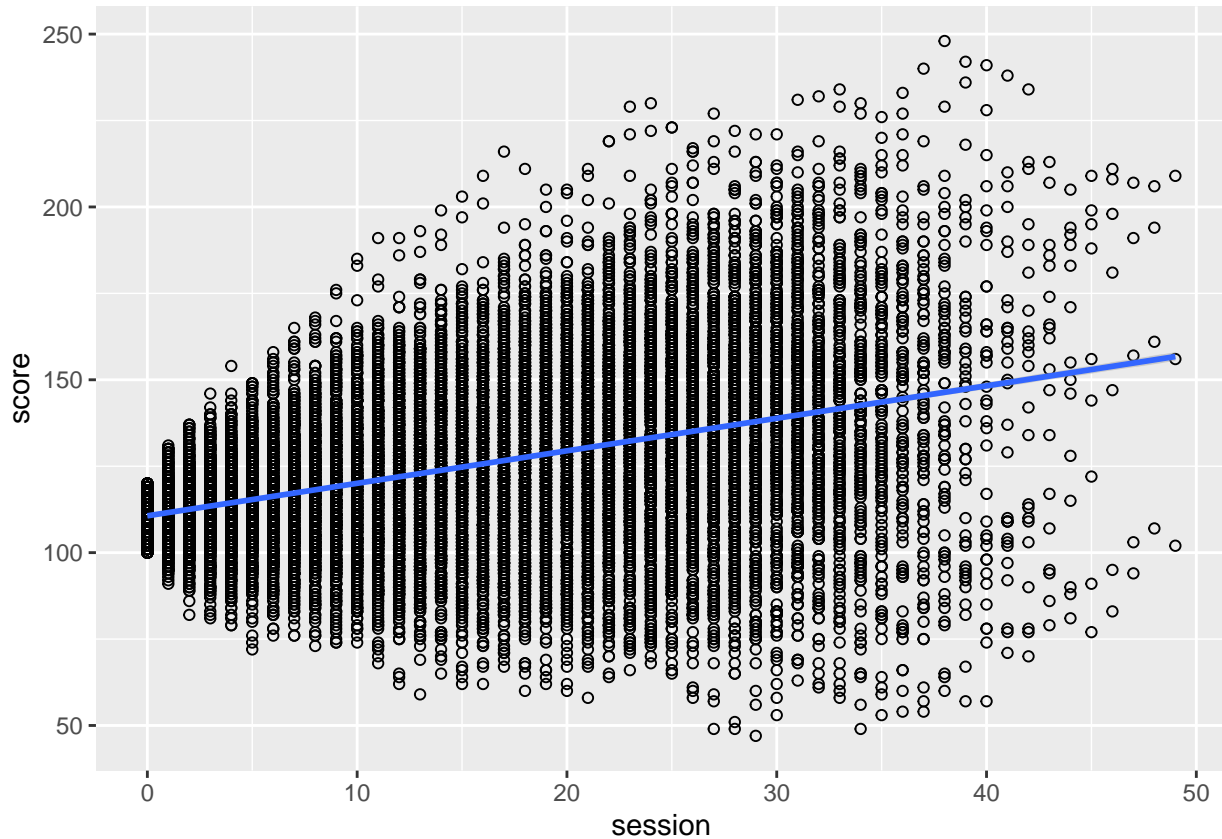
The relationship between the score and the session can be observed with the next two figures. The regression line (blue) in the scatterplot clearly shows how the score rises with the amount of sessions. This can also be observed in the box plot when looking at the mean (black line) for every box.

```
# set labels
ds$sessionF <- factor(ds$session, levels=c(0:49), labels=c(0:49))

# boxplot score per session
boxplot(score~sessionF, data=ds, main="Score", xlab="session", ylab="score")
```



```
# ggplot score per session  
hp <- ggplot(ds, aes(x=session, y=score)) + geom_point(shape=1) +  
  geom_smooth(formula = y ~ x, method=lm)  
hp
```



3.2 Frequentist approach

3.2.1 Multilevel analysis

We have conducted a multilevel analysis. We have an intercept only model (model0) which we compare to a model that includes a predictor parameter for the session (model1). By comparing these two models, we will know whether there is a difference in the score over the sessions.

```
# create models as given in slides lecture 4
model0 <- lm(formula=score~1, data=ds, na.action=na.exclude)
model1 <- lm(formula=score~sessionF, data=ds, na.action=na.exclude)

# analysis, see if predictor improves fitting
pander(anova(model0,model1))
```

Table 9: Analysis of Variance Table

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
16127	10713477	NA	NA	NA	NA
16078	9228641	49	1484836	52.79	0

```
pander(anova(model1))
```

Table 10: Analysis of Variance Table

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
sessionF	49	1484836	30303	52.79	0
Residuals	16078	9228641	574	NA	NA

From this analysis we can see there is a significant variation between the sessions. We take a further look at the summary results.

```
pander(summary(model1))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	110.3	1.07	103	0
sessionF1	0.7166	1.514	0.4734	0.636
sessionF2	1.908	1.514	1.261	0.2075
sessionF3	2.96	1.514	1.955	0.05054
sessionF4	3.838	1.514	2.536	0.01123
sessionF5	5.004	1.514	3.306	0.0009494
sessionF6	5.972	1.514	3.945	8.005e-05
sessionF7	7.016	1.514	4.635	3.599e-06
sessionF8	7.637	1.514	5.045	4.585e-07
sessionF9	8.551	1.514	5.649	1.642e-08
sessionF10	9.097	1.515	6.004	1.969e-09
sessionF11	10.36	1.515	6.836	8.453e-12
sessionF12	11.22	1.515	7.402	1.41e-13
sessionF13	12.44	1.515	8.209	2.407e-16
sessionF14	13.6	1.515	8.974	3.161e-19
sessionF15	14.45	1.515	9.539	1.644e-21
sessionF16	15.63	1.515	10.31	7.277e-25
sessionF17	16.69	1.516	11.01	4.212e-28
sessionF18	18.07	1.518	11.9	1.62e-32
sessionF19	19.18	1.521	12.61	2.751e-36
sessionF20	19.8	1.521	13.02	1.539e-38
sessionF21	20.85	1.525	13.67	2.489e-42
sessionF22	21.35	1.529	13.96	5.131e-44
sessionF23	22.39	1.536	14.58	7.593e-48
sessionF24	23.32	1.542	15.12	2.555e-51
sessionF25	25.06	1.555	16.11	5.923e-58
sessionF26	26.11	1.574	16.59	2.818e-61
sessionF27	26.54	1.597	16.62	1.578e-61
sessionF28	27.37	1.64	16.69	5.045e-62
sessionF29	28.76	1.68	17.12	4.195e-65
sessionF30	29.52	1.73	17.07	9.765e-65
sessionF31	30.66	1.804	16.99	3.552e-64
sessionF32	30.02	1.908	15.73	2.406e-55
sessionF33	30.19	2.034	14.85	1.568e-49
sessionF34	30.08	2.213	13.59	7.575e-42
sessionF35	30.84	2.388	12.92	5.712e-38
sessionF36	31.17	2.571	12.12	1.116e-33
sessionF37	29.23	2.825	10.35	5.185e-25
sessionF38	32.42	3.076	10.54	6.942e-26
sessionF39	31.77	3.767	8.434	3.618e-17
sessionF40	32.78	4.031	8.131	4.553e-16

	Estimate	Std. Error	t value	Pr(> t)
sessionF41	32.06	4.503	7.119	1.13e-12
sessionF42	34.55	5.109	6.763	1.397e-11
sessionF43	37.34	5.748	6.496	8.474e-11
sessionF44	38.8	6.492	5.976	2.33e-09
sessionF45	43.17	8.057	5.358	8.536e-08
sessionF46	50.16	9.118	5.5	3.846e-08
sessionF47	40.13	10.77	3.727	0.0001948
sessionF48	56.73	12.03	4.717	2.417e-06
sessionF49	45.39	13.87	3.272	0.00107

Table 12: Fitting linear model: score ~ sessionF

Observations	Residual Std. Error	R^2	Adjusted R^2
16128	23.96	0.1386	0.136

The summary compares the first session (intercept) with the other sessions. Looking at the estimates, we can see that compared to the first session, the scores are higher every later session. Next, we will take a look at the Akaike Information Criterion (AIC) to compare the models on the goodness-of-fit concerning the out-of-sample deviance.

```
#AIC
models <- list(model0, model1)
model.names <- c("model0", "model1")
pander(aictab(cand.set = models, modnames=model.names),
       caption="Model selection based on AICc.")
```

Table 13: Model selection based on AICc.

	Modnames	K	AICc	Delta_AICc	ModelLik	AICcWt	LL	Cum.Wt
2	model1	51	148277	0	1	1	-74087	1
1	model0	2	150584	2308	0	0	-75290	1

Here we can see that model 1 has the best goodness-of-fit as it has the smallest AICc value. Lastly, we will obtain a 95% confidence interval of the estimates we have obtained earlier.

```
# gives CI95%
pander(confint(model1), caption="95% confidence interval of the estimates.")
```

Table 14: 95% confidence interval of the estimates.

	2.5 %	97.5 %
(Intercept)	108.2	112.4
sessionF1	-2.251	3.684
sessionF2	-1.059	4.875
sessionF3	-0.007004	5.927
sessionF4	0.8712	6.805
sessionF5	2.037	7.971
sessionF6	3.005	8.939
sessionF7	4.049	9.983

	2.5 %	97.5 %
sessionF8	4.67	10.6
sessionF9	5.584	11.52
sessionF10	6.127	12.07
sessionF11	7.388	13.33
sessionF12	8.245	14.19
sessionF13	9.468	15.41
sessionF14	10.63	16.57
sessionF15	11.48	17.42
sessionF16	12.66	18.6
sessionF17	13.72	19.67
sessionF18	15.09	21.04
sessionF19	16.2	22.16
sessionF20	16.82	22.79
sessionF21	17.86	23.84
sessionF22	18.35	24.34
sessionF23	19.38	25.41
sessionF24	20.3	26.34
sessionF25	22.01	28.11
sessionF26	23.02	29.19
sessionF27	23.41	29.67
sessionF28	24.15	30.58
sessionF29	25.47	32.06
sessionF30	26.13	32.91
sessionF31	27.12	34.19
sessionF32	26.28	33.76
sessionF33	26.2	34.18
sessionF34	25.74	34.42
sessionF35	26.16	35.52
sessionF36	26.13	36.21
sessionF37	23.69	34.76
sessionF38	26.39	38.45
sessionF39	24.39	39.16
sessionF40	24.88	40.68
sessionF41	23.23	40.89
sessionF42	24.54	44.57
sessionF43	26.07	48.6
sessionF44	26.07	51.52
sessionF45	27.38	58.96
sessionF46	32.28	68.03
sessionF47	19.02	61.23
sessionF48	33.15	80.3
sessionF49	18.2	72.59

Here we can again see an increase in score related to the sessions. From this we can conclude that the session has a positive effect on people's score. Also, it seems there is a significant variance between the participants in their score when the sessions increase.

3.2.2 Report section for a scientific publication

A Linear Model analysis was conducted to test the difference between sessions on the score. The results found a significant effect ($F(49,16078) = 52.79$, $p < .001$) for the sessions on the score. From the results we can conclude that over sessions the score per participant generally increases. Moreover, the variance between

the participants increases when the sessions increase, which we think is caused by missing scores on later sessions or the ground truth. We cannot make any conclusions on this until the missing data is available.

3.3 Bayesian approach

3.3.1 Model description

For model 2, the model with session as a factor, we take as prior a normal distribution of $N(125,30)$. This comes from the mean of the score, 125, and a bit more than the standard deviation, which is around 27. Our sigma is set at a uniform distribution of $U(0.001,30)$.

$$score \sim Norm(\mu, \sigma)$$

$$\mu = a + b * sessionF + c * subjF$$

$$alpha = Norm(125, 30)$$

$$\sigma = Uniform(0.001, 50)$$

3.3.2 Model comparison

```
ds <- ds[!(ds$Subject>99),] # select first 100 subjects
ds$Subject <- ds$Subject +1 # increase subject number with 1 to overcome Stan zero index problem

mean(ds$score) # check mean

## [1] 125.5142

sd(ds$score) # check standard deviation

## [1] 27.402

ds$sessionF <- factor(ds$session, levels=c(0:49), labels=c(0:49))
ds$subjF <- factor(ds$Subject, levels=c(1:100), labels=c(1:100))

da <- subset(ds, select=c(score, sessionF))
da1 <- subset(ds, select=c(score, sessionF, subjF))

# create model with fixed intercept (i)
m0 <- map2stan(
  alist(
    score ~ dnorm(mu, sigma),
    mu <- a,
    a ~ dnorm(125,30), # mean and sd from what we found above
    sigma ~ dunif(0.001,50)
  ), data = da, iter = 10000, chains = 4, cores = 4
)

## Computing WAIC

# create model extended with an adaptive prior for subject id (ii)
m1 <- map2stan(
  alist(
```

```

score ~ dnorm(mu, sigma),
mu <- a + a_subj[subjF],
a_subj[subjF] ~ dnorm(0, sigma_subj),
sigma_subj ~ dcauchy(0, 10),
a ~ dnorm(125, 30),
sigma ~ dcauchy(0.001, 50)
), data = da1, iter = 10000, chains = 4, cores = 4
)

```

Computing WAIC

```

# create model with session as a factor (iii)
m2 <- map2stan(
  alist(
    score ~ dnorm(mu, sigma),
    mu <- a[sessionF],
    a[sessionF] ~ dnorm(125, 30),
    sigma ~ dunif(0.001, 50)
  ), data = da, iter = 10000, chains = 4, cores = 4
)

```

Computing WAIC

```
pander(compare(m0, m1, m2, func=WAIC))
```

	WAIC	SE	dWAIC	dSE	pWAIC	weight
m1	28322	96.41	0	NA	93.07	1
m2	30578	98.61	2256	120.9	67.55	0
m0	31039	98.23	2717	104.5	2.48	0

We have created and compared the three described models. From the results we can see that model 1, the model with the adaptive prior for subject id, has the best fit since it has the smallest WAIC value and largest Akaike weight.

3.3.3 Estimates examination

From the previous comparison we could see that model 1 is the best fit model. We will further examine this model with 95% credible intervals of the parameters of this model.

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

	mean	sd	2.5%	97.5%	n_eff	Rhat4
a_subj[1]	-6.652	5.84	-18.2	4.833	6141	1
a_subj[2]	-19.51	3.644	-26.53	-12.38	2339	1.001
a_subj[3]	44.41	3.695	37.23	51.78	2511	1.001
a_subj[4]	14.16	3.559	7.197	21.11	2193	1.001
a_subj[5]	-5.908	3.634	-13.08	1.181	2339	1.001
a_subj[6]	-10.66	3.824	-18.22	-3.14	2533	1.001
a_subj[7]	-4.571	3.74	-11.91	2.811	2546	1.001
a_subj[8]	9.373	3.717	2.017	16.65	2374	1.001
a_subj[9]	-2.592	4.543	-11.52	6.282	3704	1
a_subj[10]	20.65	3.768	13.25	27.99	2499	1.001
a_subj[11]	13.2	3.303	6.841	19.72	1888	1.001
a_subj[12]	-19.72	3.683	-27	-12.5	2418	1.001

	mean	sd	2.5%	97.5%	n_eff	Rhat4
a_subj[13]	9.494	3.513	2.625	16.49	2184	1.001
a_subj[14]	10.25	3.403	3.534	16.97	2055	1.001
a_subj[15]	-25.67	3.514	-32.58	-18.78	2213	1.001
a_subj[16]	26.57	3.545	19.59	33.52	2251	1.001
a_subj[17]	2.373	3.513	-4.403	9.372	2147	1.001
a_subj[18]	-42.47	3.893	-50.07	-34.88	2630	1.001
a_subj[19]	-1.155	3.868	-8.697	6.449	2591	1.001
a_subj[20]	-2.435	3.883	-10.04	5.166	2653	1.001
a_subj[21]	-13.31	3.52	-20.23	-6.304	2297	1.001
a_subj[22]	10.85	4.367	2.125	19.39	3351	1.001
a_subj[23]	-6.585	3.504	-13.45	0.1838	2114	1.001
a_subj[24]	-6.251	4.174	-14.3	1.962	2971	1
a_subj[25]	-15.55	3.782	-22.98	-8.063	2470	1.001
a_subj[26]	-7.636	3.383	-14.14	-0.9418	2077	1.001
a_subj[27]	2.23	3.492	-4.589	9.12	2276	1
a_subj[28]	7.842	3.772	0.3836	15.29	2574	1.001
a_subj[29]	-30.12	3.686	-37.26	-22.89	2360	1.001
a_subj[30]	23.82	3.738	16.58	31.32	2463	1.001
a_subj[31]	23.46	3.977	15.69	31.28	2622	1.001
a_subj[32]	-43.19	3.623	-50.28	-36.03	2408	1.001
a_subj[33]	34.21	3.549	27.27	41.22	2146	1.001
a_subj[34]	-29	3.815	-36.45	-21.39	2515	1.001
a_subj[35]	4.879	3.899	-2.828	12.52	2707	1.001
a_subj[36]	-32.12	3.364	-38.79	-25.46	2000	1.001
a_subj[37]	-4.808	3.729	-12.07	2.43	2423	1.001
a_subj[38]	20.26	3.428	13.52	26.88	2086	1.001
a_subj[39]	-12.6	3.766	-20	-5.233	2420	1.001
a_subj[40]	-10.81	3.845	-18.33	-3.332	2622	1.001
a_subj[41]	-24.6	3.857	-32.08	-16.92	2635	1
a_subj[42]	14.1	3.822	6.696	21.58	2587	1.001
a_subj[43]	-2.646	3.776	-10.05	4.846	2364	1.001
a_subj[44]	-29.46	3.808	-36.88	-21.99	2524	1.001
a_subj[45]	-17.01	3.585	-24.04	-9.978	2248	1.001
a_subj[46]	1.408	3.434	-5.249	8.217	2037	1.001
a_subj[47]	9.14	3.631	2.064	16.25	2257	1.001
a_subj[48]	-10.67	3.652	-17.85	-3.387	2361	1.001
a_subj[49]	38.72	3.539	31.82	45.63	2224	1.001
a_subj[50]	19.31	3.588	12.28	26.36	2190	1.001
a_subj[51]	1.256	3.676	-5.905	8.557	2349	1.001
a_subj[52]	-8.073	3.487	-14.85	-1.151	2184	1.001
a_subj[53]	17.04	3.748	9.778	24.4	2438	1.001
a_subj[54]	7.917	3.691	0.6416	15.1	2395	1.001
a_subj[55]	-13.58	3.997	-21.38	-5.669	2726	1.001
a_subj[56]	-29	4.033	-36.86	-21.12	2904	1.001
a_subj[57]	10.51	3.675	3.193	17.69	2416	1.001
a_subj[58]	14.18	3.786	6.804	21.68	2508	1.001
a_subj[59]	-24.94	4.16	-33.11	-16.81	2977	1.001
a_subj[60]	21.32	3.508	14.44	28.31	2080	1.001
a_subj[61]	-17.93	3.521	-24.87	-11.02	2159	1.001
a_subj[62]	72.82	3.366	66.28	79.44	2046	1.001
a_subj[63]	-34.14	3.757	-41.58	-26.8	2475	1.001
a_subj[64]	-19.61	3.845	-27.19	-12.03	2488	1.001

	mean	sd	2.5%	97.5%	n_eff	Rhat4
a_subj[65]	-27.78	3.573	-34.87	-20.79	2225	1.001
a_subj[66]	-20.4	3.462	-27.25	-13.69	2023	1.001
a_subj[67]	17.1	3.871	9.485	24.56	2771	1.001
a_subj[68]	-15.65	3.571	-22.61	-8.558	2214	1.001
a_subj[69]	13.34	3.774	5.93	20.66	2479	1.001
a_subj[70]	33.58	3.645	26.38	40.64	2229	1.001
a_subj[71]	-27.78	3.713	-34.99	-20.52	2515	1.001
a_subj[72]	-16.5	3.685	-23.68	-9.194	2413	1.001
a_subj[73]	5.965	3.856	-1.616	13.65	2630	1.001
a_subj[74]	13.35	3.641	6.262	20.54	2278	1.001
a_subj[75]	-17.89	3.325	-24.4	-11.36	1956	1.001
a_subj[76]	12.9	3.586	5.988	20.01	2244	1
a_subj[77]	29.64	3.502	22.8	36.45	2191	1.001
a_subj[78]	-5.521	3.683	-12.85	1.631	2404	1.001
a_subj[79]	-16	3.808	-23.42	-8.411	2616	1.001
a_subj[80]	-3.562	3.534	-10.56	3.397	2153	1.001
a_subj[81]	6.879	3.665	-0.319	14.03	2282	1.001
a_subj[82]	-9.766	3.579	-16.77	-2.783	2222	1.001
a_subj[83]	-3.99	3.89	-11.52	3.767	2653	1.001
a_subj[84]	-33.39	3.619	-40.54	-26.43	2306	1.001
a_subj[85]	-7.883	3.762	-15.33	-0.4665	2481	1.001
a_subj[86]	20.39	3.409	13.8	27.13	1984	1.001
a_subj[87]	7.464	3.824	-0.07451	14.92	2660	1.001
a_subj[88]	21.25	3.931	13.5	28.86	2710	1.001
a_subj[89]	-1.285	3.474	-7.989	5.538	2136	1.001
a_subj[90]	19.13	3.758	11.78	26.49	2454	1.001
a_subj[91]	-0.8107	3.66	-8.015	6.384	2454	1.001
a_subj[92]	15.89	3.804	8.543	23.39	2656	1.001
a_subj[93]	3.441	3.636	-3.695	10.57	2401	1.001
a_subj[94]	3.012	3.473	-3.678	9.863	2242	1.001
a_subj[95]	-8.377	4.157	-16.4	-0.2763	3101	1.001
a_subj[96]	1.086	3.867	-6.469	8.686	2713	1.001
a_subj[97]	29.27	3.93	21.56	36.98	2646	1.001
a_subj[98]	24.26	3.466	17.5	31.09	2122	1.001
a_subj[99]	0.2742	3.548	-6.628	7.226	2260	1.001
a_subj[100]	14.62	3.571	7.597	21.58	2250	1.001
sigma_subj	20.49	1.51	17.76	23.69	25637	1
a	125	2.051	120.9	128.9	763.9	1.003
sigma	17.86	0.2247	17.43	18.31	29752	0.9999

We can observe that the mean between the subjects has a high variance. This means that although the scores increase per session, the subjects have very different prior skills, achieving relatively higher or lower scores in their first session than average.