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

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

20/04/2021

## Contents

## 1   Introduction

We will be using the following packages:

```
library(ggplot2) # plotting
library(AICcmodavg) # aictab
library(pander) #for rendering output
library(rethinking) # for stan
```

```
## Loading required package: rstan

## Loading required package: StanHeaders

## rstan (Version 2.21.2, GitRev: 2e1f913d3ca3)

## For execution on a local, multicore CPU with excess RAM we recommend calling
## options(mc.cores = parallel::detectCores()).
```

```
## To avoid recompilation of unchanged Stan programs, we recommend calling
## rstan_options(auto_write = TRUE)

## Do not specify '-march=native' in 'LOCAL_CPPFLAGS' or a Makevars file

## Loading required package: parallel

## Loading required package: dagitty

## rethinking (Version 2.01)

##
## Attaching package: 'rethinking'

## The following object is masked from 'package:AICcmodavg':
##
##     DIC

## The following object is masked from 'package:stats':
##
##     rstudent
```

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

## 2.1 The motivation for the planned research

(Max 250 words)

## 2.2 The theory underlying the research

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

## 2.3 Research questions

The research question that will be examined in the experiment (or alternatively the hypothesis that will be tested in the experiment)

## 2.4 The related conceptual model

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

## 2.5 Experimental Design

Note that the study should have a true experimental design

## 2.6 Experimental procedure

Describe how the experiment will be executed step by step

## 2.7 Measures

Describe the measure that will be used

## 2.8 Participants

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

## 2.9 Suggested statistical analyses

Describe the statistical test you suggest to care out on the collected data
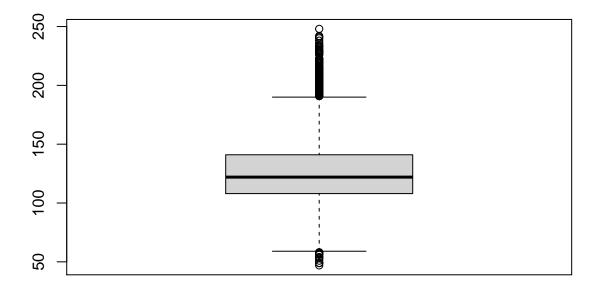
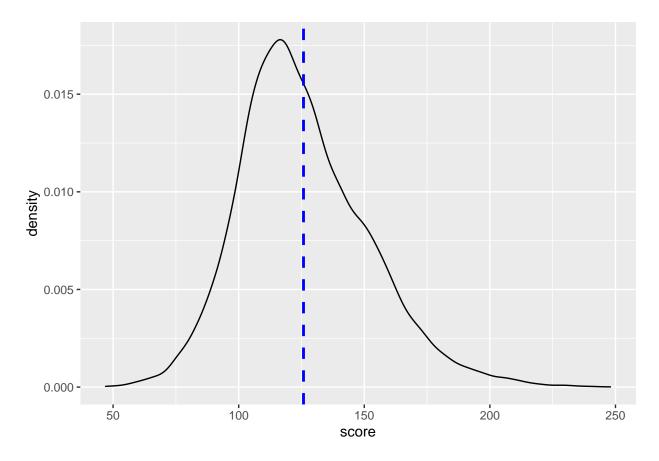# 3 Part 3 - Multilevel model

## 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)
```
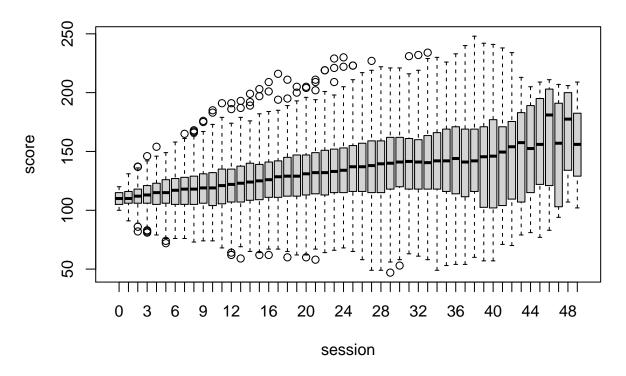


```
# density score overall distribution (with mean line)
p <- ggplot(ds, aes(x=score)) + geom_density()
p + geom_vline(aes(xintercept=mean(score)), color="blue", linetype="dashed",size=1)
```
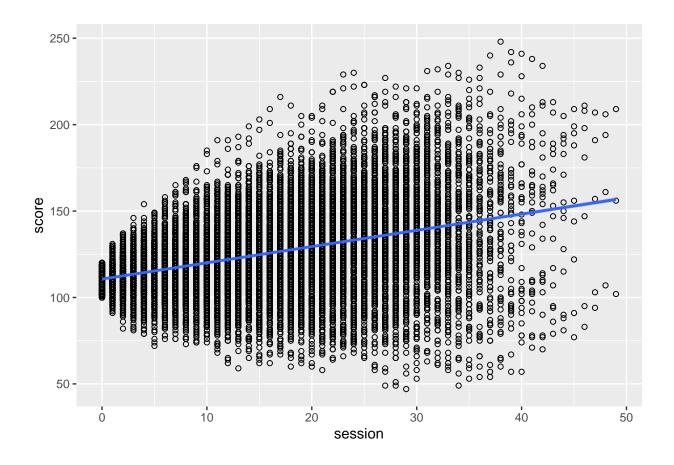
The relationship between the score and the session can be observed with the next two figures. The regression line (blue) in the scatterplot clearly shows how the score rises with the amount of sessions. This can also be observed in the box plot when looking at the mean (black line) for every box.

```
# 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")
```

**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 1: 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 2: 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 4: 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 concering 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 5: 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 6: 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 ($F_{(49,16078)} = 52.793$, $p < .001$) 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.

## 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).

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

$$\mu = \alpha$$

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

$$\sigma = Uniform(0.001, 30)$$

### 3.3.2 Model comparison

We will create and compare 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.

```
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(0.001,30)
  ), 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(0.001,30)
  ), 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(0.001,30)
  ), 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 | 96.43 | 0     | NA    | 93.39  | 1      |
| **m2** | 30579 | 98.63 | 2256  | 121   | 68.13  | 0      |
| **m0** | 31039 | 98.24 | 2716  | 104.5 | 2.448  | 0      |

### 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]**  | -6.669 | 5.807 | -17.97 | 4.719  | 7529  | 1     |
| **a_subj[2]**  | -19.54 | 3.602 | -26.54 | -12.5  | 2911  | 1.001 |
| **a_subj[3]**  | 44.45  | 3.673 | 37.31  | 51.65  | 2978  | 1.001 |
| **a_subj[4]**  | 14.12  | 3.547 | 7.106  | 21.03  | 3009  | 1     |
| **a_subj[5]**  | -5.942 | 3.635 | -13.09 | 1.161  | 2926  | 1.001 |
| **a_subj[6]**  | -10.65 | 3.76  | -18.05 | -3.357 | 3164  | 1.001 |
| **a_subj[7]**  | -4.599 | 3.766 | -11.92 | 2.98   | 3062  | 1.001 |
| **a_subj[8]**  | 9.376  | 3.665 | 2.194  | 16.55  | 3091  | 1.001 |
| **a_subj[9]**  | -2.595 | 4.527 | -11.53 | 6.205  | 4777  | 1     |
| **a_subj[10]** | 20.64  | 3.716 | 13.39  | 27.98  | 3242  | 1.001 |
| **a_subj[11]** | 13.17  | 3.316 | 6.764  | 19.72  | 2552  | 1.001 |
| **a_subj[12]** | -19.74 | 3.697 | -27    | -12.57 | 3090  | 1.001 |
| **a_subj[13]** | 9.482  | 3.533 | 2.617  | 16.45  | 2755  | 1.001 |

|  | mean | sd | 2.5% | 97.5% | n_eff | Rhat4 |
|---|---|---|---|---|---|---|
| **a_subj[14]** | 10.24 | 3.371 | 3.581 | 16.77 | 2622 | 1.001 |
| **a_subj[15]** | -25.68 | 3.477 | -32.39 | -18.79 | 2770 | 1.001 |
| **a_subj[16]** | 26.55 | 3.615 | 19.49 | 33.64 | 2774 | 1.001 |
| **a_subj[17]** | 2.372 | 3.422 | -4.284 | 9.097 | 2646 | 1.001 |
| **a_subj[18]** | -42.45 | 3.911 | -50.1 | -34.69 | 3392 | 1.001 |
| **a_subj[19]** | -1.171 | 3.812 | -8.473 | 6.354 | 3297 | 1.001 |
| **a_subj[20]** | -2.495 | 3.879 | -10.07 | 5.076 | 3422 | 1.001 |
| **a_subj[21]** | -13.32 | 3.501 | -20.11 | -6.378 | 2756 | 1.001 |
| **a_subj[22]** | 10.86 | 4.364 | 2.335 | 19.44 | 4386 | 1.001 |
| **a_subj[23]** | -6.628 | 3.503 | -13.44 | 0.2604 | 2799 | 1.001 |
| **a_subj[24]** | -6.234 | 4.176 | -14.43 | 1.879 | 3931 | 1.001 |
| **a_subj[25]** | -15.57 | 3.837 | -23.19 | -8.126 | 3309 | 1.001 |
| **a_subj[26]** | -7.63 | 3.35 | -14.26 | -1.131 | 2439 | 1.001 |
| **a_subj[27]** | 2.196 | 3.507 | -4.733 | 9.063 | 2769 | 1.001 |
| **a_subj[28]** | 7.827 | 3.699 | 0.534 | 15.05 | 3157 | 1.001 |
| **a_subj[29]** | -30.11 | 3.634 | -37.25 | -22.95 | 3105 | 1.001 |
| **a_subj[30]** | 23.78 | 3.742 | 16.44 | 31.12 | 3382 | 1.001 |
| **a_subj[31]** | 23.42 | 3.946 | 15.72 | 31.2 | 3467 | 1.001 |
| **a_subj[32]** | -43.21 | 3.613 | -50.29 | -36.21 | 3085 | 1 |
| **a_subj[33]** | 34.19 | 3.458 | 27.36 | 40.91 | 2742 | 1.001 |
| **a_subj[34]** | -29.04 | 3.813 | -36.46 | -21.55 | 3438 | 1.001 |
| **a_subj[35]** | 4.905 | 3.829 | -2.562 | 12.46 | 3262 | 1.001 |
| **a_subj[36]** | -32.17 | 3.351 | -38.8 | -25.57 | 2548 | 1.001 |
| **a_subj[37]** | -4.831 | 3.693 | -12.01 | 2.457 | 3030 | 1.001 |
| **a_subj[38]** | 20.28 | 3.425 | 13.72 | 27.17 | 2598 | 1.001 |
| **a_subj[39]** | -12.6 | 3.777 | -20.01 | -5.205 | 3202 | 1.001 |
| **a_subj[40]** | -10.89 | 3.859 | -18.44 | -3.339 | 3329 | 1 |
| **a_subj[41]** | -24.6 | 3.816 | -31.99 | -17.01 | 3302 | 1 |
| **a_subj[42]** | 14.09 | 3.811 | 6.773 | 21.67 | 3108 | 1.001 |
| **a_subj[43]** | -2.652 | 3.739 | -9.992 | 4.662 | 3281 | 1.001 |
| **a_subj[44]** | -29.48 | 3.77 | -36.77 | -22.08 | 3251 | 1.001 |
| **a_subj[45]** | -17.06 | 3.544 | -24.05 | -10.21 | 2769 | 1.001 |
| **a_subj[46]** | 1.359 | 3.434 | -5.357 | 8.022 | 2686 | 1.001 |
| **a_subj[47]** | 9.135 | 3.589 | 2.157 | 16.18 | 2849 | 1.001 |
| **a_subj[48]** | -10.73 | 3.56 | -17.75 | -3.849 | 3022 | 1.001 |
| **a_subj[49]** | 38.7 | 3.524 | 31.83 | 45.66 | 2778 | 1.001 |
| **a_subj[50]** | 19.28 | 3.58 | 12.22 | 26.34 | 2848 | 1.001 |
| **a_subj[51]** | 1.231 | 3.677 | -5.892 | 8.46 | 3099 | 1.001 |
| **a_subj[52]** | -8.091 | 3.443 | -14.9 | -1.439 | 2751 | 1.001 |
| **a_subj[53]** | 16.98 | 3.735 | 9.632 | 24.26 | 3197 | 1.001 |
| **a_subj[54]** | 7.898 | 3.612 | 0.7928 | 14.98 | 2947 | 1.001 |
| **a_subj[55]** | -13.6 | 3.988 | -21.4 | -5.619 | 3748 | 1 |
| **a_subj[56]** | -29.03 | 4.004 | -36.97 | -21.25 | 3754 | 1 |
| **a_subj[57]** | 10.5 | 3.678 | 3.384 | 17.67 | 3052 | 1.001 |
| **a_subj[58]** | 14.17 | 3.755 | 6.863 | 21.52 | 3109 | 1.001 |
| **a_subj[59]** | -24.96 | 4.16 | -33.12 | -16.86 | 3546 | 1.001 |
| **a_subj[60]** | 21.28 | 3.519 | 14.39 | 28.16 | 2747 | 1.001 |
| **a_subj[61]** | -17.97 | 3.53 | -24.96 | -11.09 | 2782 | 1.001 |
| **a_subj[62]** | 72.82 | 3.336 | 66.22 | 79.34 | 2564 | 1.001 |
| **a_subj[63]** | -34.17 | 3.718 | -41.44 | -26.8 | 3166 | 1.001 |
| **a_subj[64]** | -19.64 | 3.804 | -27.01 | -12.14 | 3164 | 1.001 |
| **a_subj[65]** | -27.83 | 3.542 | -34.73 | -20.89 | 2919 | 1.001 |

|  | mean | sd | 2.5% | 97.5% | n_eff | Rhat4 |
|---|---|---|---|---|---|---|
| a_subj[66] | -20.4 | 3.429 | -27.16 | -13.69 | 2724 | 1.001 |
| a_subj[67] | 17.04 | 3.82 | 9.55 | 24.57 | 3411 | 1 |
| a_subj[68] | -15.68 | 3.483 | -22.52 | -8.865 | 2658 | 1.001 |
| a_subj[69] | 13.3 | 3.701 | 6.059 | 20.58 | 3181 | 1.001 |
| a_subj[70] | 33.58 | 3.602 | 26.58 | 40.65 | 2910 | 1.001 |
| a_subj[71] | -27.84 | 3.682 | -34.98 | -20.72 | 2964 | 1.001 |
| a_subj[72] | -16.55 | 3.615 | -23.69 | -9.428 | 3021 | 1.001 |
| a_subj[73] | 5.944 | 3.872 | -1.721 | 13.52 | 3496 | 1 |
| a_subj[74] | 13.31 | 3.569 | 6.165 | 20.32 | 2980 | 1.001 |
| a_subj[75] | -17.93 | 3.266 | -24.36 | -11.48 | 2470 | 1.001 |
| a_subj[76] | 12.87 | 3.532 | 5.978 | 19.82 | 3021 | 1.001 |
| a_subj[77] | 29.62 | 3.45 | 22.89 | 36.39 | 2569 | 1.001 |
| a_subj[78] | -5.54 | 3.645 | -12.66 | 1.635 | 3056 | 1 |
| a_subj[79] | -16.02 | 3.781 | -23.37 | -8.615 | 3362 | 1.001 |
| a_subj[80] | -3.572 | 3.515 | -10.42 | 3.4 | 2791 | 1.001 |
| a_subj[81] | 6.833 | 3.577 | -0.1393 | 13.96 | 3030 | 1.001 |
| a_subj[82] | -9.78 | 3.526 | -16.7 | -2.855 | 2941 | 1.001 |
| a_subj[83] | -3.994 | 3.889 | -11.62 | 3.555 | 3289 | 1.001 |
| a_subj[84] | -33.41 | 3.599 | -40.42 | -26.25 | 2990 | 1.001 |
| a_subj[85] | -7.912 | 3.765 | -15.24 | -0.5427 | 3120 | 1 |
| a_subj[86] | 20.38 | 3.335 | 13.88 | 26.91 | 2394 | 1.001 |
| a_subj[87] | 7.456 | 3.75 | 0.1394 | 14.77 | 3400 | 1 |
| a_subj[88] | 21.23 | 3.87 | 13.63 | 28.78 | 3516 | 1.001 |
| a_subj[89] | -1.317 | 3.444 | -8.109 | 5.492 | 2700 | 1.001 |
| a_subj[90] | 19.12 | 3.711 | 11.85 | 26.5 | 3220 | 1.001 |
| a_subj[91] | -0.8308 | 3.65 | -7.93 | 6.381 | 3046 | 1.001 |
| a_subj[92] | 15.87 | 3.758 | 8.488 | 23.14 | 3141 | 1.001 |
| a_subj[93] | 3.41 | 3.63 | -3.703 | 10.57 | 2838 | 1.001 |
| a_subj[94] | 2.956 | 3.49 | -3.933 | 9.795 | 2698 | 1.001 |
| a_subj[95] | -8.419 | 4.132 | -16.54 | -0.228 | 3837 | 1.001 |
| a_subj[96] | 1.09 | 3.876 | -6.544 | 8.708 | 3675 | 1 |
| a_subj[97] | 29.24 | 3.863 | 21.54 | 36.73 | 3504 | 1.001 |
| a_subj[98] | 24.25 | 3.441 | 17.49 | 30.95 | 2669 | 1.001 |
| a_subj[99] | 0.2547 | 3.477 | -6.516 | 7.069 | 2898 | 1.001 |
| a_subj[100] | 14.6 | 3.462 | 7.861 | 21.47 | 2785 | 1.001 |
| sigma_subj | 20.5 | 1.501 | 17.81 | 23.7 | 29042 | 1 |
| a | 125 | 1.993 | 121 | 128.9 | 917.1 | 1.003 |
| sigma | 17.86 | 0.2255 | 17.43 | 18.31 | 36272 | 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.