

# Preliminary Analysis

## Bias in PERSUADE

```
library(sjPlot)
library(lme4)
library(car) # level names in contrasts
library(stargazer) # LaTeX tables
library(emmeans)
library(performance) # ICC
library(tidyverse)

df <- read.csv("../data/persuade_corpus.csv",
               na.strings = c("", " ", "NA")) %>%
  select(-c(full_text, X, assignment, source_text)) %>%
  mutate_if(is.character, as.factor) %>%

  # drop American Indian/Alaskan Native group because there are < 200 samples
  filter(race_ethnicity != "American Indian/Alaskan Native") %>%

  # remove unused levels (American Indian)
  droplevels() %>%

  # simplify level and variable names
  mutate(
    race=fct_recode(race_ethnicity,
                    asian="Asian/Pacific Islander",
                    black="Black/African American",
                    hisp="Hispanic/Latino",
                    other="Two or more races/Other",
                    white="White")
  ) %>%

  # collapse NaNs into negative level of binary factors
```

```

mutate(
  is_ell=fct_collapse(addNA(ell), Yes = "Yes", No = c("No", NA)),
  is_disadvantaged=fct_collapse(addNA(economically_disadvantaged),
                                Yes = "Economically disadvantaged",
                                No = c("Not economically disadvantaged", NA)),
  has_disability=fct_collapse(addNA(student_disability_status),
                              Yes = "Identified as having disability",
                              No = c("Not identified as having disability", NA))
) %>%

# set the first level, which will be used as "reference"
mutate(
  is_disadvantaged=fct_relevel(is_disadvantaged, "No"),
  has_disability=fct_relevel(has_disability, "No"),
)

# configure contrasts for race and source
options(decorate.contr.Sum = c("", ""))
contrasts(df$race) = contr.Sum(levels(df$race))
contrasts(df$source) = contr.Sum(levels(df$source))

summary(df)

```

essay_id		holistic_score_1	holistic_score_2
5.04604E+12:	2	Min. :1.000	Min. :1.0
5.88194E+12:	2	1st Qu.:2.000	1st Qu.:2.0
2021000039 :	1	Median :3.000	Median :3.0
2021000047 :	1	Mean :3.272	Mean :3.3
2021000071 :	1	3rd Qu.:4.000	3rd Qu.:4.0
2021000080 :	1	Max. :6.000	Max. :6.0
(Other)	:25847		

holistic_score_adjudicated	source
Min. :1.000	Florida :3987
1st Qu.:2.000	Georgia Virtual:1165
Median :3.000	Indiana :8799
Mean :3.317	NCES :4798
3rd Qu.:4.000	Virginia :7106
Max. :6.000	

	prompt_name	task	gender
Distance learning	: 2153	Independent :13069	F:13067

Facial action coding system : 2150    Text dependent:12786    M:12788  
 Does the electoral college work?: 2035  
 Car-free cities : 1952  
 Driverless cars : 1868  
 Exploring Venus : 1849  
 (Other) :13848

	grade	ell	race_ethnicity
Min.	: 6.000	No :22318	Asian/Pacific Islander : 1743
1st Qu.:	8.000	Yes : 2242	Black/African American : 4959
Median :	9.000	NA's: 1295	Hispanic/Latino : 6560
Mean :	9.174		Two or more races/Other: 1022
3rd Qu.:	10.000		White :11571
Max.	:12.000		

economically\_disadvantaged  
 Economically disadvantaged : 9565  
 Not economically disadvantaged:11074  
 NA's : 5216

	student_disability_status	race	is_ell
Identified as having disability	: 3325	asian: 1743	No :23613
Not identified as having disability:	21365	black: 4959	Yes: 2242
NA's	: 1165	hisp : 6560	
		other: 1022	
		white:11571	

is_disadvantaged	has_disability
No :16290	No :22530
Yes: 9565	Yes: 3325

```

eval_model <- function(mod) {
  print(performance(mod))
}
  
```

```
# print(tab_model(mod,
#                 p.adjust="HB",
#                 show.aic=TRUE,
#                 show.re.var=TRUE,
#                 # show.reflvl=TRUE,
#                 prefix.labels="varname"
#                 ))
}
```

```
mod.null <-lmer(holistic_score_adjudicated ~ 1 + (1|prompt_name), data=df)
eval_model(mod.null)
```

# Indices of model performance

AIC	AICc	BIC	R2 (cond.)	R2 (marg.)	ICC	RMSE	Sigma
73488.393	73488.394	73512.874	0.270	0.000	0.270	1.000	1.000

## Simple model with just race

```
mod.race <-lmer(holistic_score_adjudicated
~ race
+ (race|prompt_name),
data=df
)
```

boundary (singular) fit: see help('isSingular')

```
eval_model(mod.race)
```

Random effect variances not available. Returned R2 does not account for random effects.

# Indices of model performance

AIC	AICc	BIC	R2 (cond.)	R2 (marg.)	RMSE	Sigma
72069.770	72069.805	72241.135		0.044	0.970	0.971

## All fixed effects

```
mod.all_fixed_effects <- lmer(holistic_score_adjudicated
  ~ race
  + grade
  + is_disadvantaged
  + is_ell
  + has_disability
  # + source # slightly improves model fit, but is not of interest
  + gender
  + (1|prompt_name),
  data=df
)
eval_model(mod.all_fixed_effects)
```

# Indices of model performance

AIC		AICc		BIC		R2 (cond.)		R2 (marg.)		ICC		RMSE		Sigma
69367.123		69367.135		69465.046		0.387		0.121		0.303		0.922		0.922

## Search for interactions

We drop grade because it was not significant, but we test for an interaction between grade and ell. We might expect interactions: race\*ses, race\*ell, and race\*disability.

```
mod.interactions <- lmer(holistic_score_adjudicated
  ~ race
  + is_disadvantaged
  + is_ell
  + has_disability
  + gender
  # + race*is_disadvantaged # increases AIC
  + race*is_ell
  # + race*has_disability # increases AIC
  # + race*source # increases AIC
  # + race*gender # increases AIC
  # + race*grade # increases AIC
  + is_disadvantaged*has_disability)
```

```

+ is_disadvantaged*is_ell
+ has_disability*is_ell
# + gender*has_disability # increases AIC
# + gender*is_disadvantaged # increases AIC
# + gender*is_ell # small decrease (-1) to AIC
+ (1|prompt_name),
  data=df
)
eval_model(mod.interactions)

```

# Indices of model performance

AIC		AICc		BIC		R2 (cond.)		R2 (marg.)		ICC		RMSE		Sigma
69272.702		69272.729		69419.587		0.394		0.125		0.308		0.920		0.920

## Search for random slopes

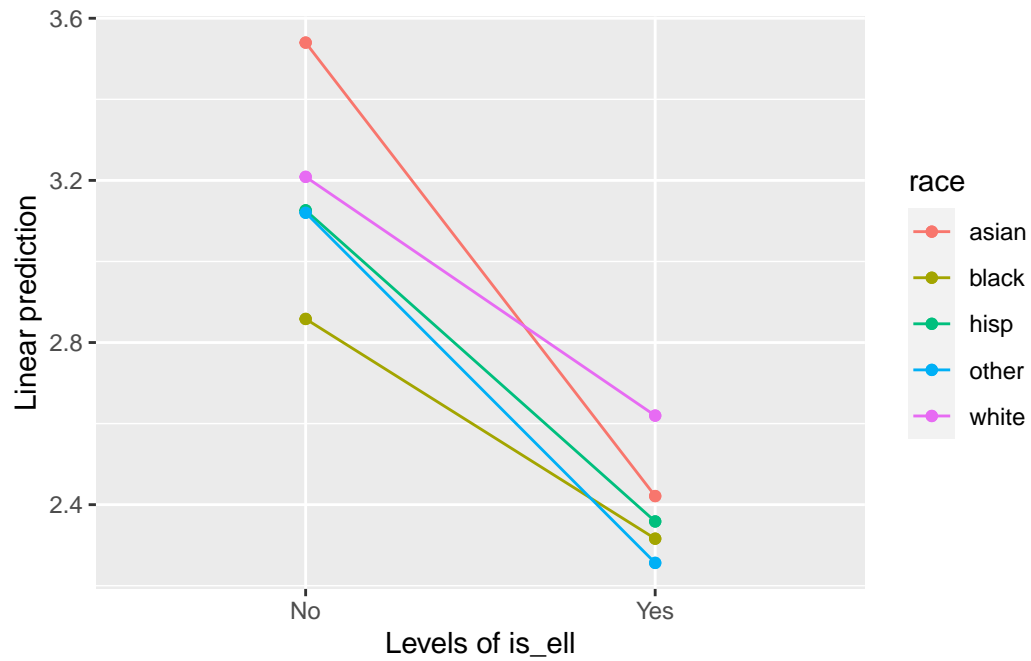
There are no good random slopes.

```

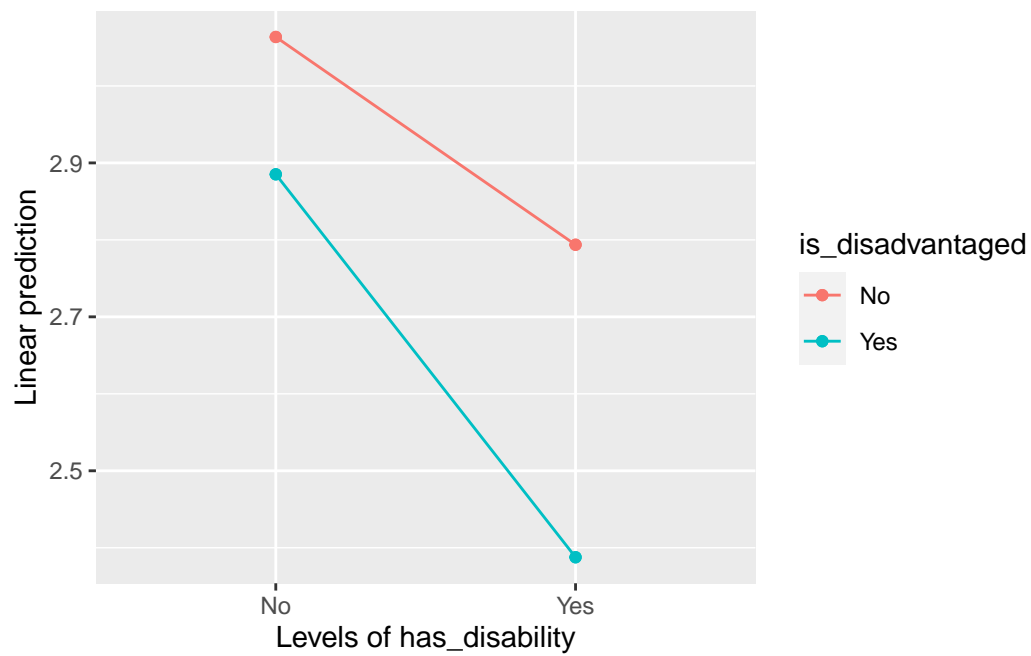
mod.final <- mod.interactions

emmip(mod.final, race~is_ell, mode = "asyp")

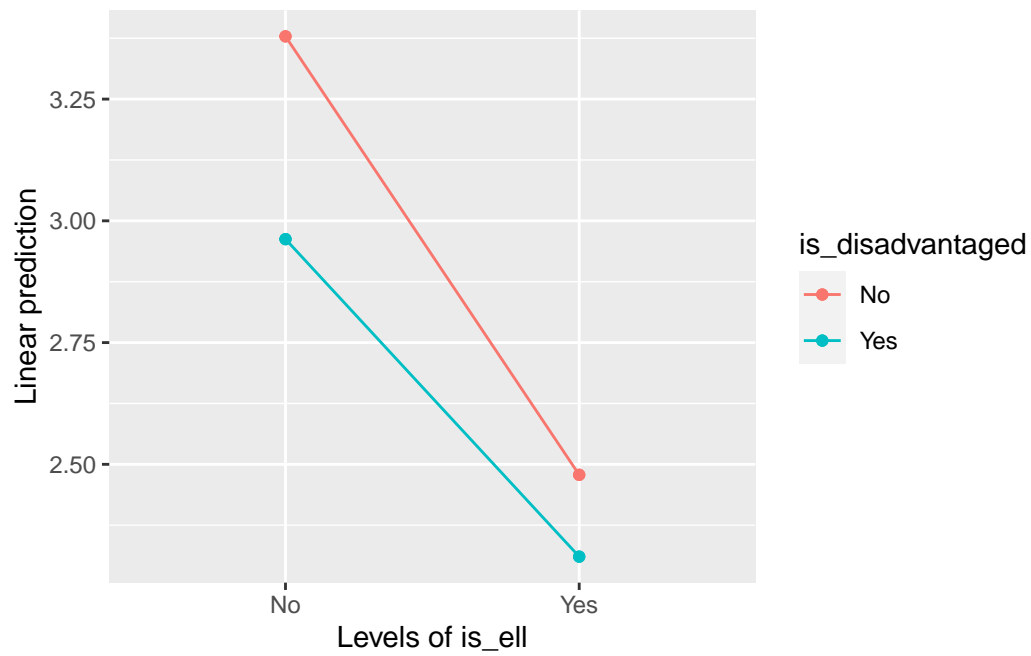
```



```
emmip(mod.final, is_disadvantaged~has_disability, mode = "asympt")
```

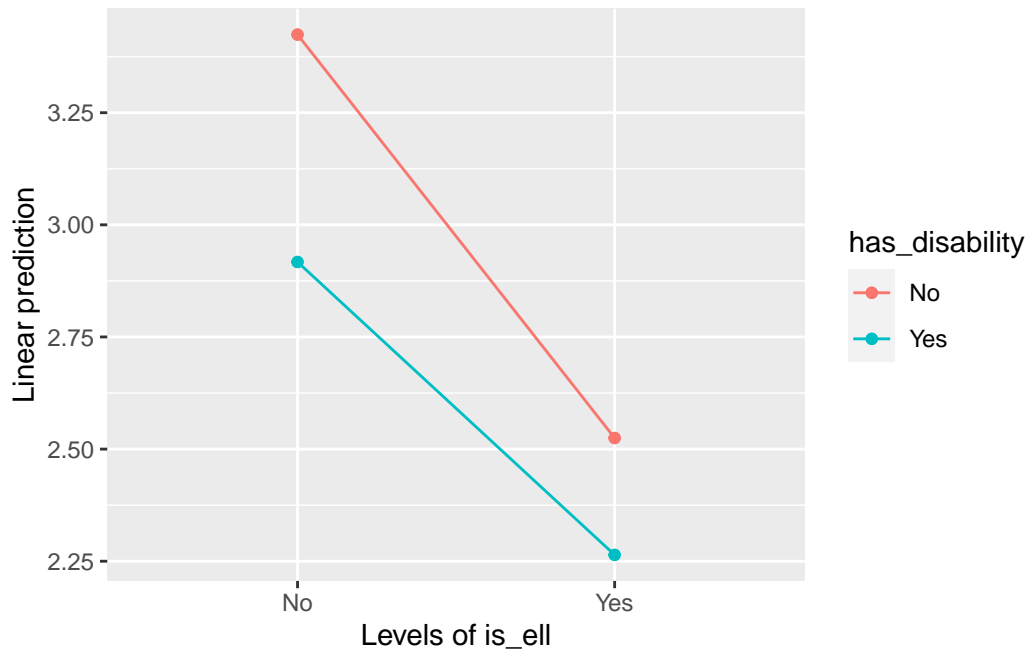


```
emmip(mod.final, is_disadvantaged~is_ell, mode = "asympt")
```



```
emmip(mod.final, has_disability~is_ell, mode = "asympt")
```





```
coefs <- summary(mod.final)$coefficients
cat(rownames(coefs), sep="\\", "\\")
```

(Intercept)", "race[asian]", "race[black]", "race[hisp]", "race[other]", "is\_disadvantagedYes"

```
labels <- c("(Intercept)", "Asian/Pacific Islander", "Black/African American", "Hispanic/L
```

```
tab_model(mod.final,
  title = "Essay Scores Regressed on Demographic Variables",
  dv.labels = "Holistic Score",
  pred.labels = labels,
  emph.p = TRUE,
  p.adjust = "BH",
  # show.reflvl=TRUE,
  show.re.var=TRUE,
  file="../results/RQ1.html"
)
```

Length of `pred.labels` does not equal number of predictors, no labelling applied.

## Essay Scores Regressed on Demographic Variables

Holistic Score

Predictors

Estimates

CI

p

(Intercept)

3.70

3.39 – 4.01

**<0.001**

race[asian]

0.37

0.33 – 0.41

**<0.001**

race[black]

-0.31

-0.34 – -0.29

**<0.001**

race[hisp]

-0.04

-0.07 – -0.02

**0.001**

race[other]

-0.05

-0.10 – -0.00

**0.041**

is\_disadvantagedYes

-0.30

-0.33 – -0.27

**<0.001**

is\_ellYes

-1.02

-1.16 – -0.89

**<0.001**

has\_disabilityYes

-0.39

-0.44 – -0.35

**<0.001**

genderM

-0.25

-0.27 – -0.22

**<0.001**

race[asian]:is\_ellYes

-0.34

-0.49 – -0.19

**<0.001**

race[black]:is\_ellYes

0.23

0.07 – 0.39

**0.005**

race[hisp]:is\_ellYes

0.01

-0.12 – 0.14

0.893

race[other]:is\_ellYes

-0.09

-0.53 – 0.35

0.743  
 is\_disadvantagedYes:has\_disabilityYes  
 -0.23  
 -0.30 – -0.16  
**<0.001**  
 is\_disadvantagedYes:is\_ellYes  
 0.25  
 0.16 – 0.34  
**<0.001**  
 is\_ellYes:has\_disabilityYes  
 0.25  
 0.13 – 0.36  
**<0.001**  
 Random Effects  
 2  
 0.85  
 00 prompt\_name  
 0.38  
 ICC  
 0.31  
 N<sub>prompt\_name</sub>  
 15  
 Observations  
 25855  
 Marginal R<sup>2</sup> / Conditional R<sup>2</sup>  
 0.125 / 0.394