

An Investigation of Grade 3-5 Students' State Test Scores in the Denver, Colorado Area

STAA 566: Final Project

Frances Lin

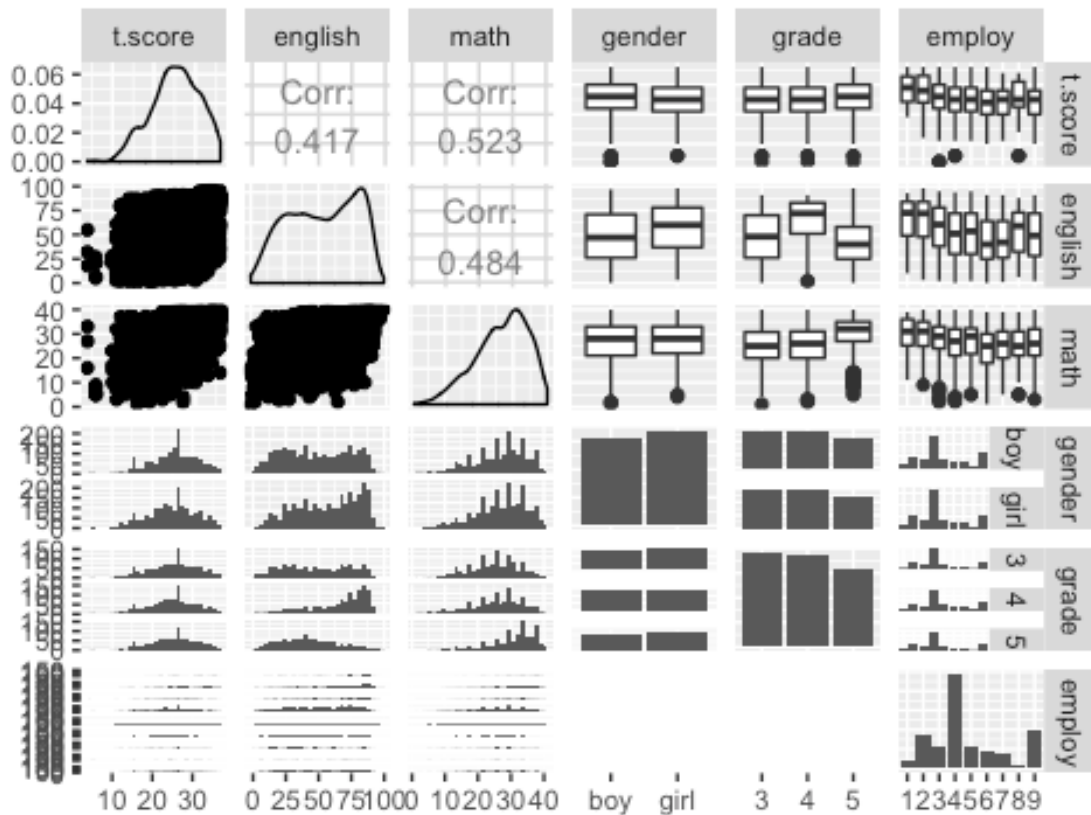
5/2019

I. DESCRIPTION

The purpose of the study is to investigate the relationship between state test score of elementary school students in the Denver, Colorado area and other variables thought to be related. Data were provided as part of the project, which included school ID, student ID, state test score, English portion of the score, math portion of the score, gender, student's grade, and parents' employment for a total of 1402 elementary school students in the Denver, Colorado area ($n=1402$). A total of 3236 test scores were recorded ($obs=3236$). State test score is the main response to test, and the rest of the variables are predictors of interest.

State test score ranges from 4 to 36 and has a mean equal to 25.13 and a median equal to 25. English portion of the score ranges from 0 to 98 and has a mean equal to 52.49 and a median equal to 54, and math portion of the score ranges from 1 to 40 and has a mean equal to 26.66 and a median equal to 28. Gender has two levels (1: boy, 2: girl). Grade ranges from 3 to 5 (1154 scores were taken in grade 3; 1129 scores were taken in grade 4; and 953 scores were taken in grade 5). Parents' employment has 9 levels (1: both parents employed, manual laborers; 2: father manual labor, mother non-manual labor; 3: father non-manual labor, mother manual labor; 4: both parents non-manual labor; 5: father employed, mother unemployed; 6: father unemployed, mother employed; 7: both parents long term unemployed; 8: both parents currently employed; 9: father absent).

Pairwise Scatterplots

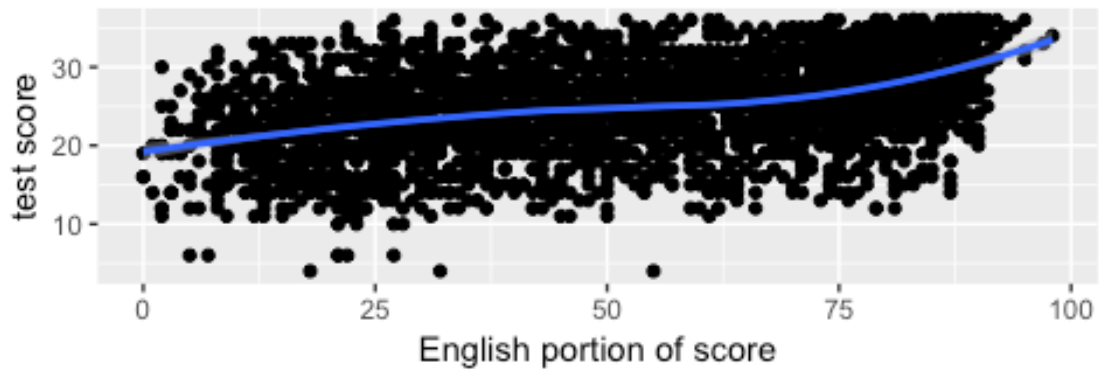


II. ANALYSIS

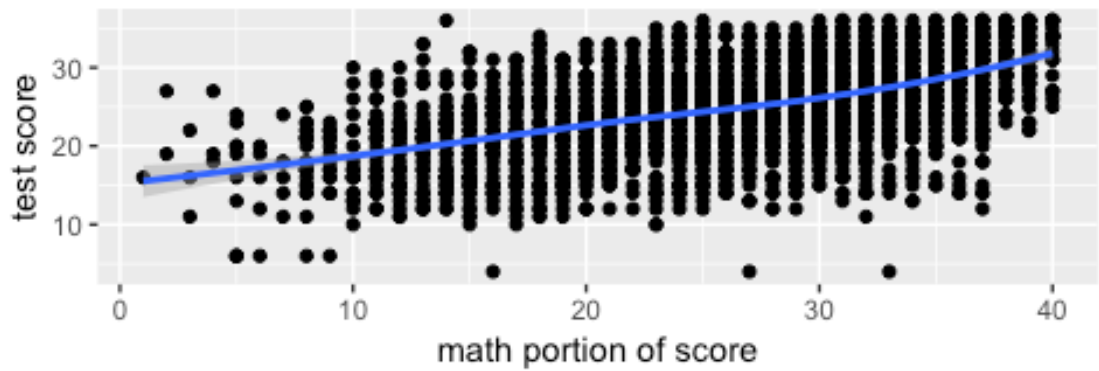
1. The first part of the analysis investigates the relationship between state test score (t.score) and English portion of the score (english) vs. the relationship between state test score (t.score) math portion of the score (math).

Initially, it appears that the relationship is consistent across test portion. Specifically, a higher English score is associated with a higher test score, a higher math score is associated with a higher test score, and both relationships appear linear. (It is worth noting that the LOESS model does not fit well, and log-transformation is recommended.) However, when parents' employment type is controlled for, it appears that the relationship is not consistent across test portion. For example, when both parents are non-manual labors (employ=4), the relationship appears consistent. On the other hand, when both parents are currently employed (employ=8), one relationship appears polynomial, and the other appears logarithmic.

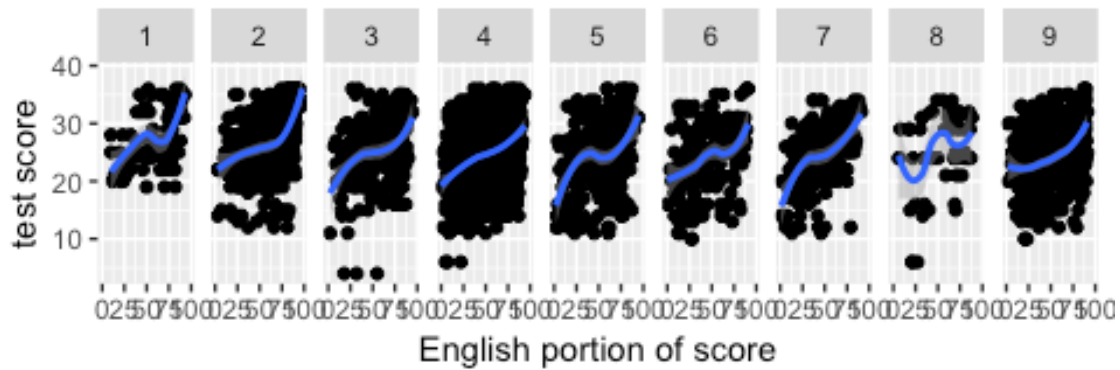
Plot of Test Score by English Score



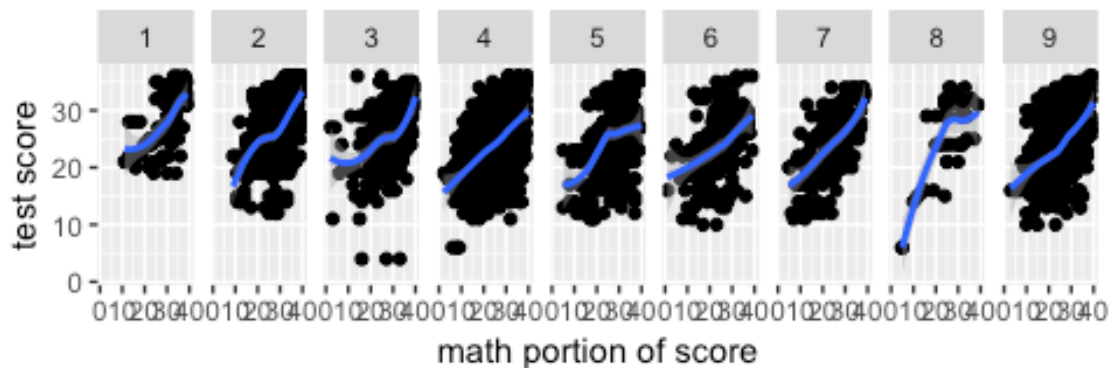
Plot of Test Score by Math Score



Plot of Test Score by English Score based on Parents' E



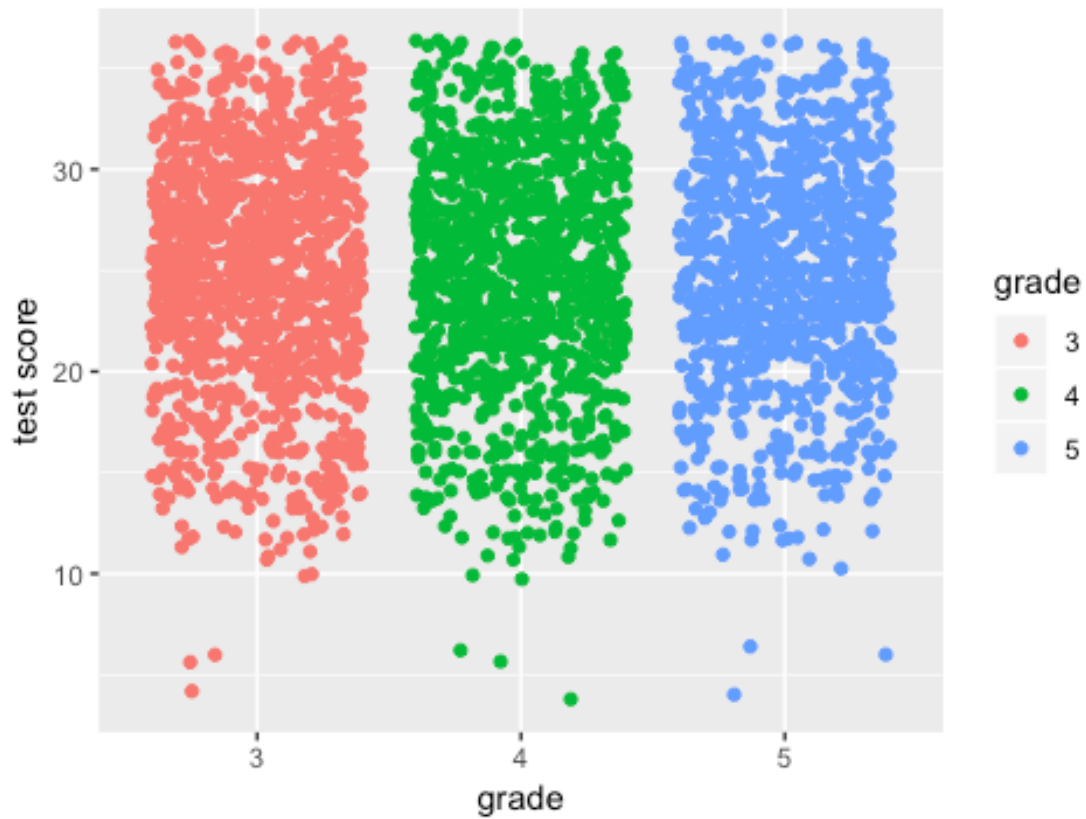
Plot of Test Score by Math Score based on Parents' Em



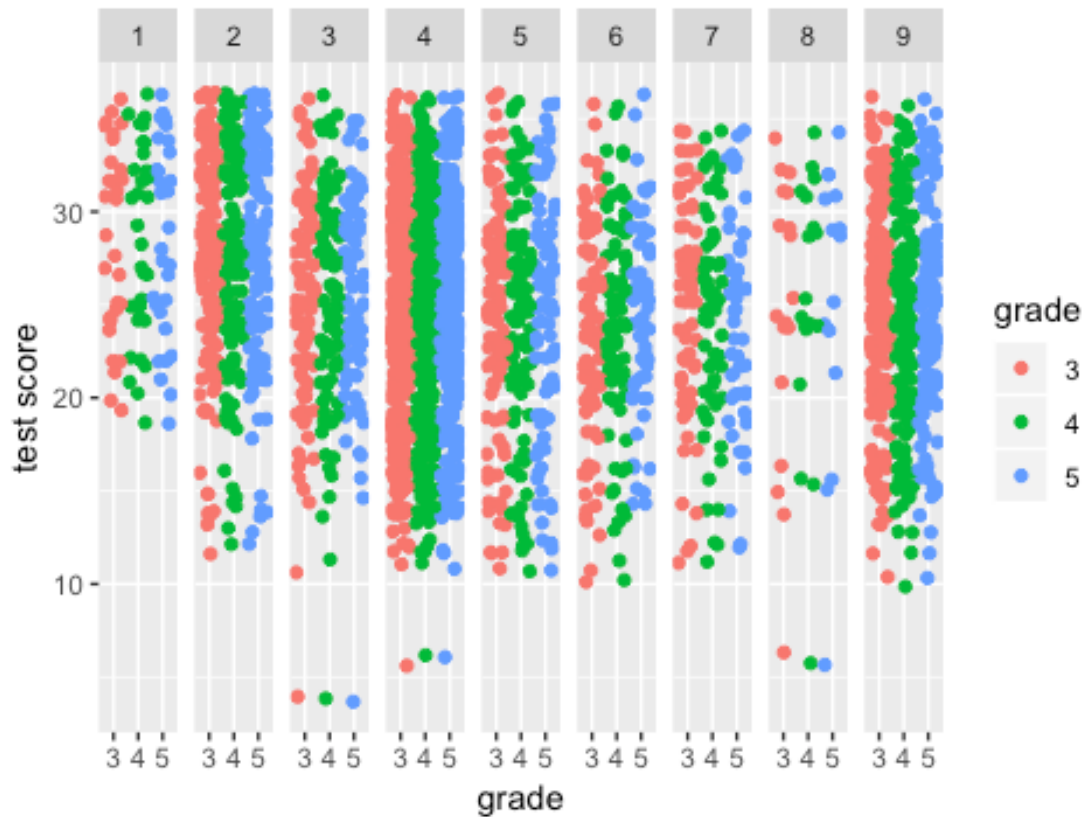
2. The second part of the analysis investigates the relationship between state test score (t.score) and student's grade (grade).

Initially, it appears that there is significant mean difference for at least one test score ($p=0.4244$). However, when parents' employment type is controlled for, it appears that there is no significant mean difference in test scores ($p=0.5039$). For example, when both parents are currently employed (employ=8), although the boxplot shows that mean test score for grade 5 is significantly higher than mean test score for grade 3 and that for grade 4, it is only because this group has sparse data.

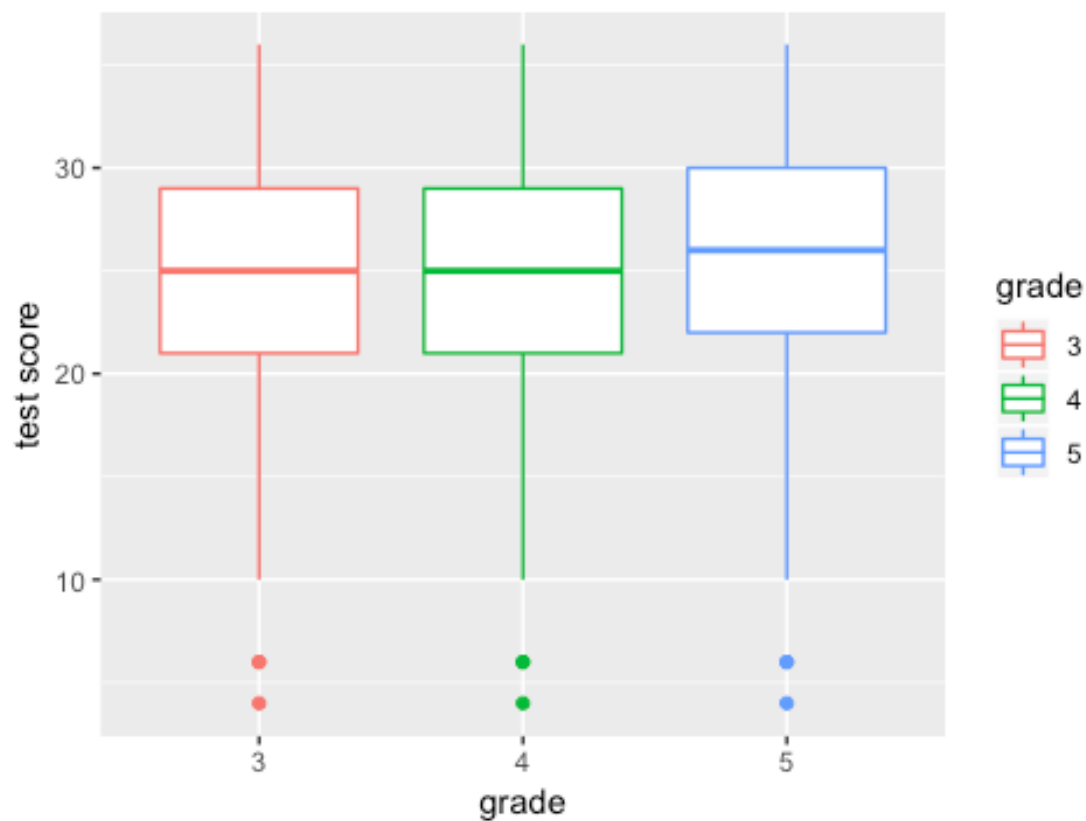
Plot of Test Score by Grade



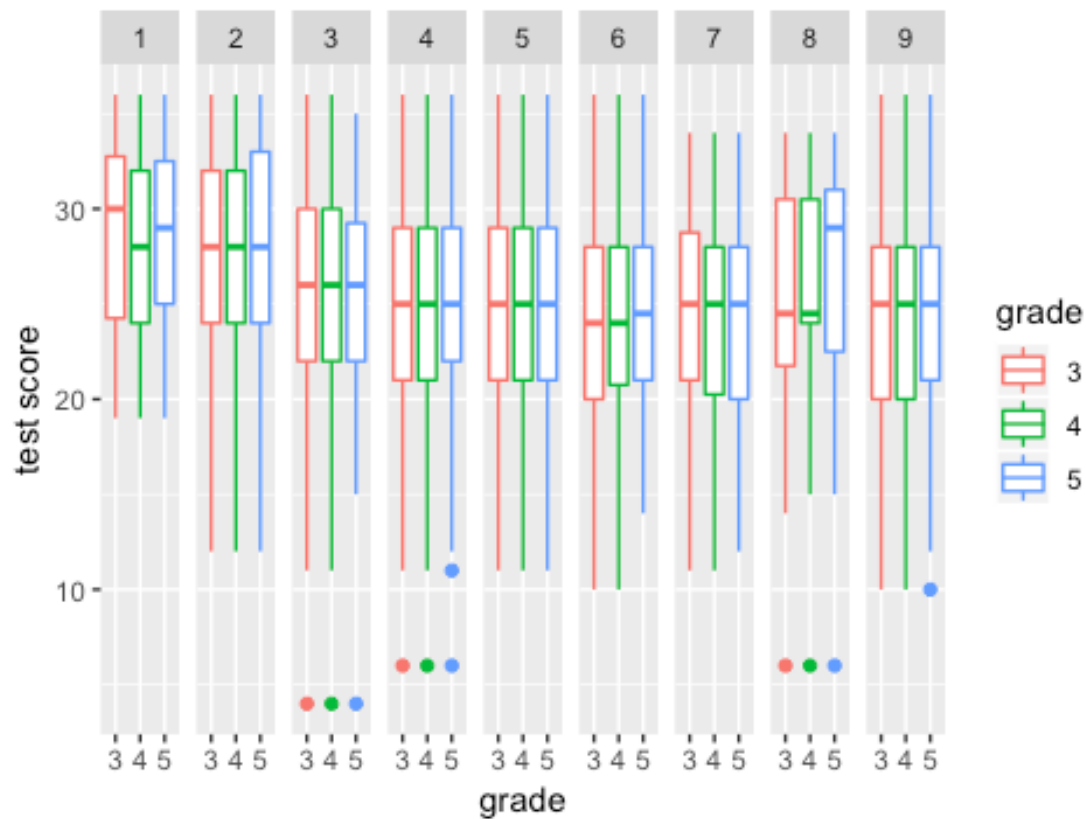
Plot of Test Score by Grade based on Parents' Employn



Plot of Test Score by Grade



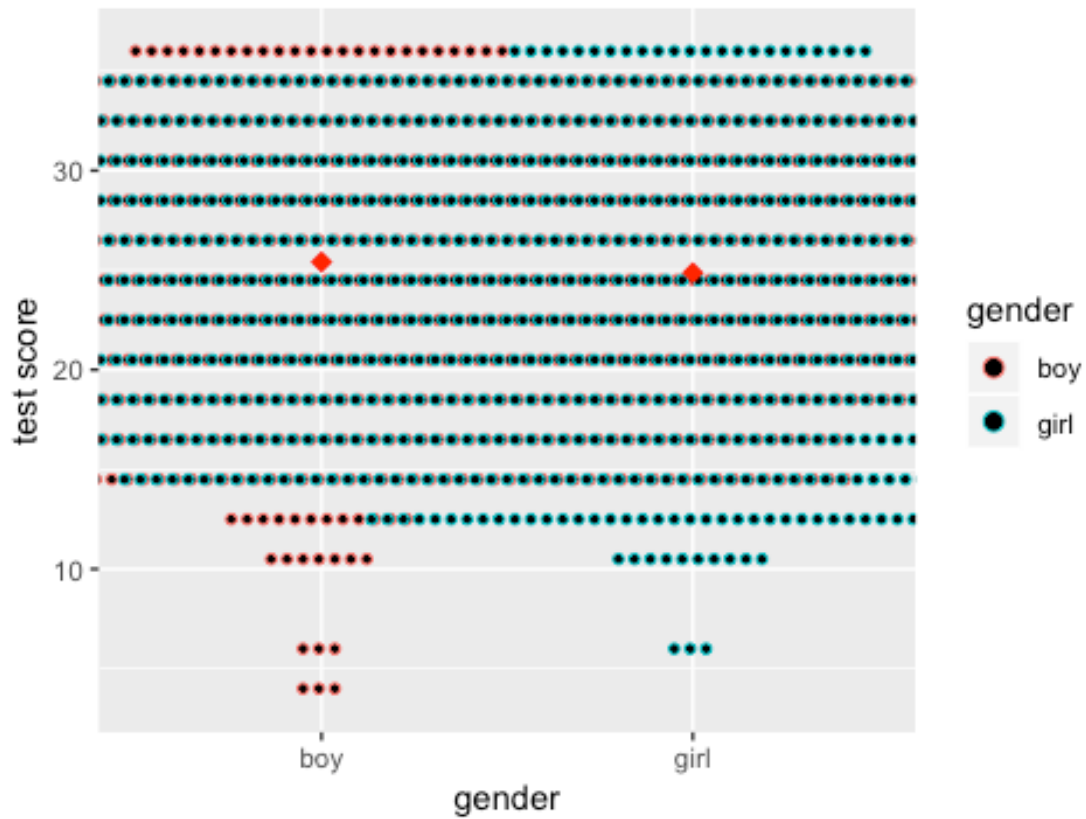
Plot of Test Score by Grade based on Parents' Employn



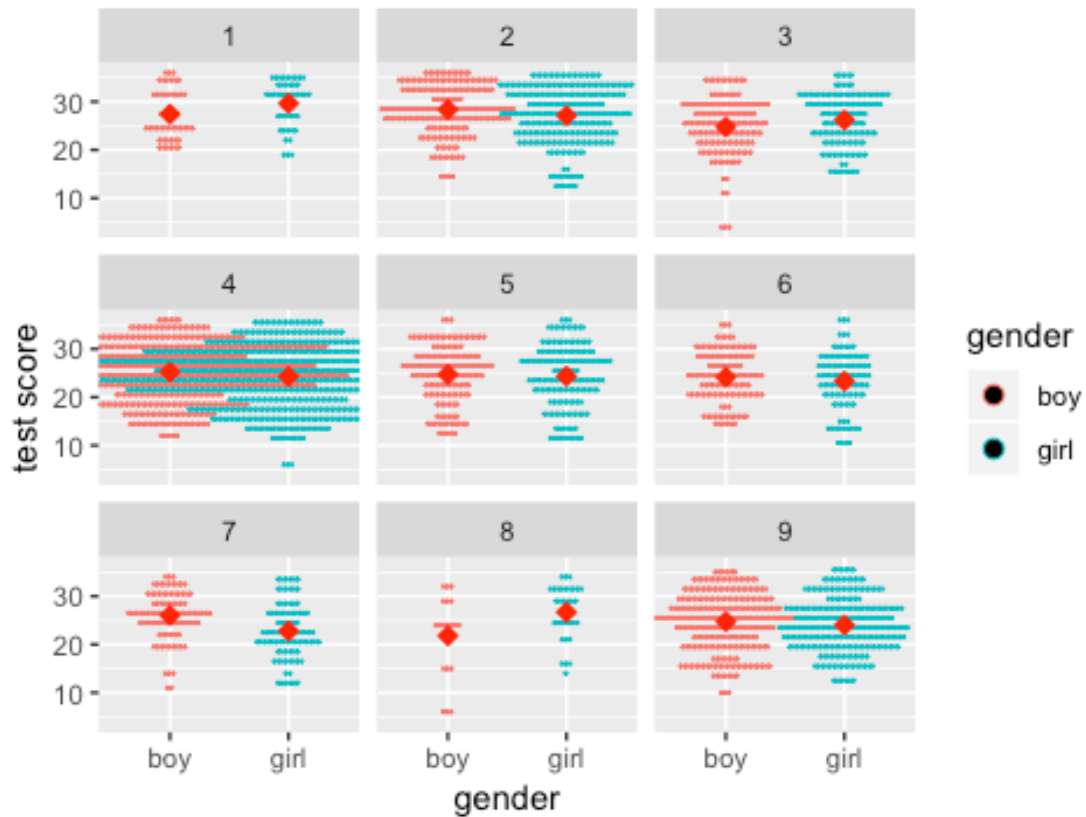
3. The third part of the analysis investigates the relationship between state test score (t.score) and student's gender (gender).

Initially, it appears that there is significant mean difference for at least one test score ($p = 0.007251$ **). When parents' employment type is controlled for, it appears that there is significant mean difference in test scores ($p = 0.001711$ **). For example, when both parents are currently employed (employ=8), the spreads are clearly different.

Plot of Test Score by Gender



Plot of Test Score by Gender based on Parents' Employ



III. RESULTS

It appears that it is not necessarily the case that the higher the English portion or the math portion of the score is, the higher the overall state test score is. In addition, it appears that parents' employment type is a confounding variable, whether the model uses grade or gender as the main predictor of interest. It is also possible that parents' employment type interacts with one or more variables. Further studies should attempt to fit separate model 1) using English portion of the score as the response to test, 2) using math portion of the score as the response to test, 3) using parents' employment type as the main predictor of interest, and 4) allowing for interaction to see if the effect of one variable on the response variable depends on the other variable, and vice versa.

IV. R CODE

```
library(GGally)
library(ggplot2)
library(gridExtra)
library(car)

# Load data
ScoreData<-read.csv(file.choose(), header=TRUE)
str(ScoreData)
View(ScoreData)

# Change variable(s) to factor
ScoreData$school.id <- as.factor(ScoreData$school.id)
ScoreData$employ <- as.factor(ScoreData$employ)
ScoreData$ID <- as.factor(ScoreData$ID)
ScoreData$grade <- as.factor(ScoreData$grade)
#str(ScoreData)

# Attach data set
attach(ScoreData)

# Summary Statistics
summary(ScoreData)

# Reorder columns
ScoreData <- ScoreData[c(1,5,4,6,7,2,8,3)]
#View(ScoreData)

# Pairwise Scatterplots using ggplot2
ggpairs(ScoreData[c(-1,-2)]) + ggtitle("Pairwise Scatterplots")

# Ques 1a: x=english, y=t.score, based on employ
g1 <- ggplot(ScoreData, aes(x=english, y=t.score)) +
  labs(title="Plot of Test Score by English Score", x="English portion
of score", y="test score")
geng0 <- g1 + geom_point() + geom_smooth(method="loess")
geng <- g1 + geom_point() + geom_smooth(method="loess") +
  facet_grid(.~employ) + ggtitle("Plot of Test Score by English Score
based on Parents' Employment")

# Ques 1b: x=math, y=t.score, based on employ
g11 <- ggplot(ScoreData, aes(x=math, y=t.score)) +
  labs(title="Plot of Test Score by Math Score", x="math portion of
score", y="test score")
gm0 <- g11 + geom_point() + geom_smooth(method="loess")
gm <- g11 + geom_point() + geom_smooth(method="loess") +
  facet_grid(.~employ) + ggtitle("Plot of Test Score by Math Score
based on Parents' Employment")

# Combine Plots of Test Score by English Score and by Math Score
grid.arrange(geng0, gm0, nrow=2)
```

```

grid.arrange(geng, gmath, nrow=2)

# Ques 2: x=grade, y=t.score, based on employ
g2 <- ggplot(ScoreData, aes(x=grade, y=t.score, color=grade)) +
  labs(title="Plot of Test Score by Grade", y="test score")
g2 + geom_jitter()
g2 + geom_jitter() + facet_grid(.~employ) +
  ggtitle("Plot of Test Score by Grade based on Parents' Employment")
g2 + geom_boxplot()
g2 + geom_boxplot() + facet_grid(.~employ) +
  ggtitle("Plot of Test Score by Grade based on Parents' Employment")

# Ques 3: x=gender, y=t.score, based on employ
g3 <- ggplot(ScoreData, aes(x=gender, y=t.score, color=gender)) +
  labs(title="Plot of Test Score by Gender", y="test score")
g3 + geom_dotplot(binaxis='y', stackdir='center', stackratio=1.5,
  dotsize=0.5) +
  stat_summary(fun.y=mean, geom="point", shape=18, size=3, color="red")
g3 + geom_dotplot(binaxis='y', stackdir='center', stackratio=1.5,
  dotsize=0.5) +
  stat_summary(fun.y=mean, geom="point", shape=18, size=3, color="red")
+
  facet_wrap(~employ) +
  ggtitle("Plot of Test Score by Gender based on Parents' Employment")

# Fit simple linear regression model for Ques 1
lmeng <- lm(t.score~english)
lmmath <- lm(t.score~math)
summary(lmeng, type=3)
summary(lmmath, type=3)

lmengc <- lm(t.score~english+employ)
lmmathc <- lm(t.score~math+employ)
Anova(lmengc, type=3)
Anova(lmmathc, type=3)

# Fit simple linear regression model for Ques 2
lmgrade <- lm(t.score~grade)
lmgradec <- lm(t.score~grade+employ)
Anova(lmgrade, type=3)
Anova(lmgradec, type=3)

# Fit simple linear regression model for Ques 3
lmgender <- lm(t.score~gender)
lmgenderc <- lm(t.score~gender+employ)
Anova(lmgender, type=3)
Anova(lmgenderc, type=3)

```