

# Project1

Yahui Li

2020/1/13

## Step1 Read Data

```
#install.packages("AER")
library(AER)
data("STAR")
```

## Step2 Explore Data

We will only examine the math scores in 1st grade in this project.

```
data <- data.frame(star1 = STAR$star1, math1 = STAR$math1)
```

```
sapply(data,class)
```

```
##      star1      math1
## "factor" "integer"
```

```
sapply(data,summary)
```

```
## $star1
##      regular      small regular+aide      NA's
##      2584      1925      2320      4769
##
## $math1
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
##      404.0  500.0  529.0  530.5  557.0  676.0  4998
```

```
data.star1.na <- data[is.na(data$star1),]
all(is.na(data.star1.na$math1))
```

```
## [1] TRUE
```

Which shows that the math score has not been recorded if class type is not recorded. So we can remove the data where star1 is NA.

One of the way to deal with NA in math1 is to remove them

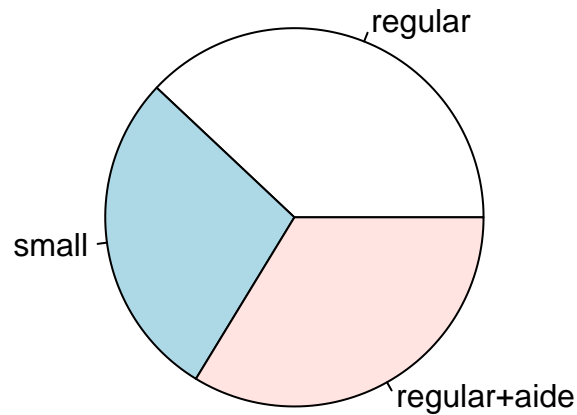
```
data_remove_na <- na.omit(data[-is.na(data$star1),])
```

```
table(data_remove_na$star1)
```

```
##
##      regular      small regular+aide
##      2507      1868      2225
```

```
pie(table(data_remove_na$star1),main = "pie chart of STAR class type")
```

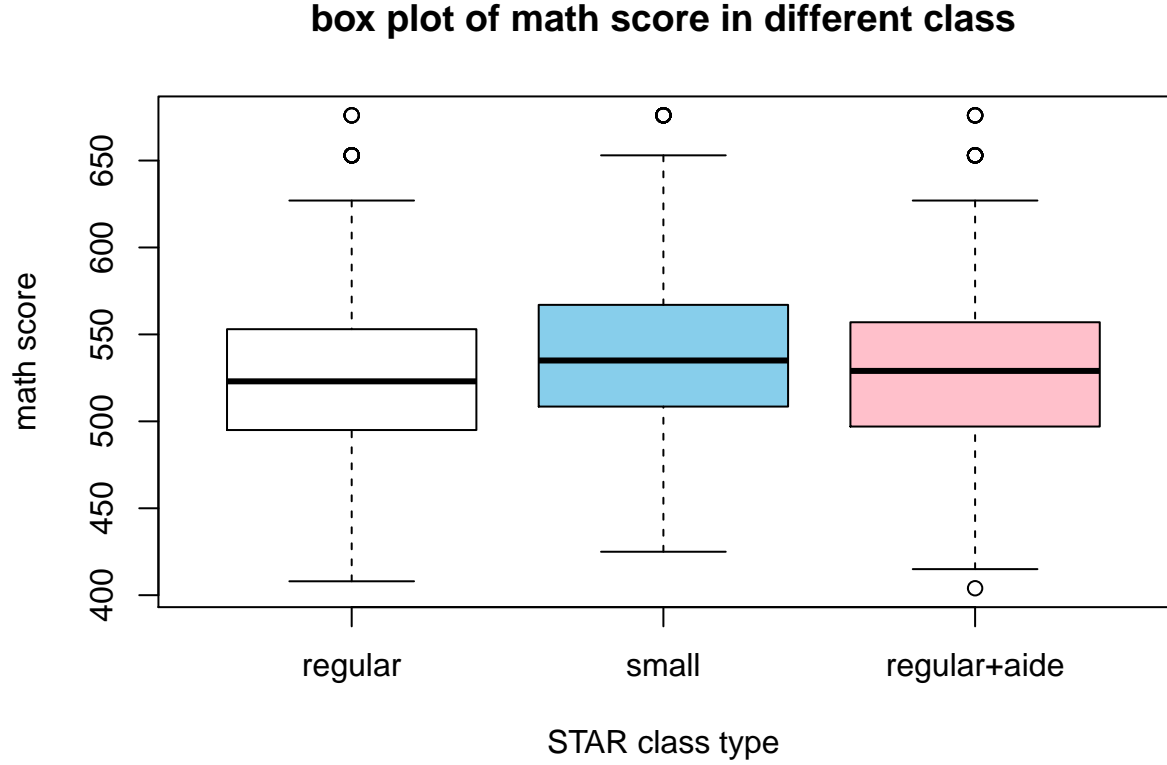
## pie chart of STAR class type



```
tapply(data_remove_na$math1, data_remove_na$star1,summary)
```

```
## $regular
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  408.0  495.0   523.0   525.3  553.0   676.0
##
## $small
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  425.0  509.2   535.0   538.7  567.0   676.0
##
## $`regular+aide`
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  404.0  497.0   529.0   529.6  557.0   676.0
```

```
boxplot(data$math1~data$star1,main = "box plot of math score in different class",
        xlab = "STAR class type", ylab = "math score", col = c("white", "skyblue", "pink"))
```



From the result,

for mean,  $\text{small} > \text{regular+aide} > \text{regular}$ ;

for all quantile information,  $\text{small} > \text{regular+aide} > \text{regular}$ ;

for min,  $\text{small} > \text{other two}$ ; For max, they are the same.

Something interesting: there are only some certain scores like 601 612 627 653 676.

### Step3 One Way ANOVA Model

$$Y_{ij} = \mu_1 + \tau_2 X_{2,ij} + \tau_3 X_{3,ij} + \epsilon_{ij}, \quad \epsilon_{ij} \sim N(0, \sigma^2), \quad i = 1, 2, 3, j = 1, \dots, n_i.$$

where  $i = 1$  means the class type in 1st grade is regular;  $i = 2$  means the class type in 1st grade is small;  $i = 3$  means the class type in 1st grade is regular-with-aide.

From the table in step2,  $n_1 = 2507, n_2 = 1868, n_3 = 2225, n = 6600$ .

$X_{2,ij} = 1$  if  $i = 2$ , otherwise  $X_{2,ij} = 0$ .  $X_{3,ij} = 1$  if  $i = 3$ , otherwise  $X_{3,ij} = 0$ .

$Y_{ij}$  denotes the math grade in 1st grade of the  $j$ -th experimental unit in the  $i$ -th class type.

$\mu_i$  means the population mean of the  $i$ -th type class in 1st grade,  $i = 1, 2, 3$ .

$\tau_i = \mu_i - \mu_1$  means the difference in population mean between  $i$ -th type and first type in 1st grade,  $i = 2, 3$ .

$\epsilon_{ij}$  is independent and identically distributed normal random variable with 0 mean and  $\sigma^2$  variance under normal assumption.

#### Model Assumption

- (a) Response variable residuals are normally distributed.
- (b) Variances of populations are equal.
- (c) Responses for a given group are independent and identically distributed normal random variables.

All of the assumptions are necessary, because for each population violate the normal distribution, F-tests are not robust. Moreover, if the assumption of homoscedasticity is violated, the Type I error properties degenerate much more severely.[5] ( Randolph, E. A.; Barcikowski, R. S. (1989). "Type I error rate when real study values are used as population parameters in a Monte Carlo study". Paper presented at the 11th annual meeting of the Mid-Western Educational Research Association, Chicago.)

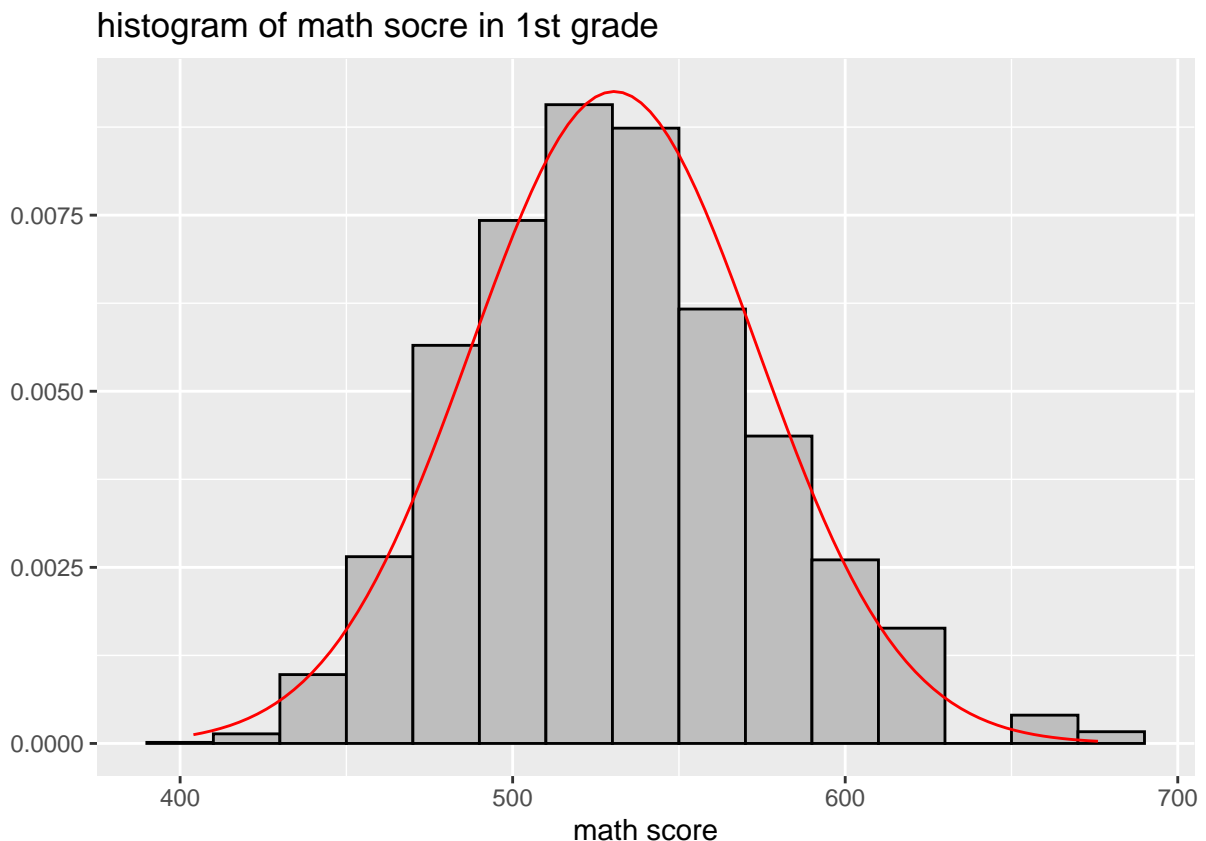
#### Step4 Appropriate

Before we fit the model, we need to ensure that model is appropriate on this dataset, that is, the response variable satisfies the assumptions of our model. In other words, we will check the normality and equal variance of the response variables.

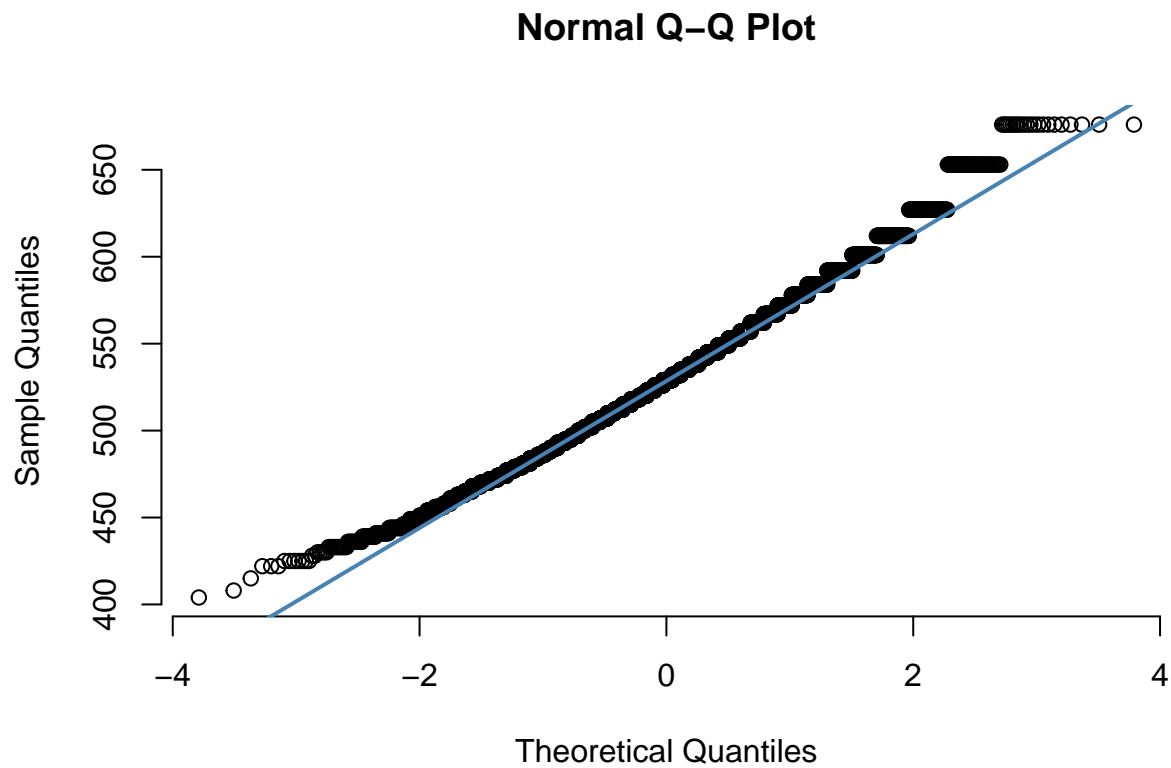
We first make a density plot and a Q-Q plot to check the normality of math1.

```
library(ggplot2)
x <- seq(404, 676, length.out=100)
df <- with(data_remove_na, data.frame(x = x, y = dnorm(x, mean(math1), sd(math1)))))

ggplot(data_remove_na, aes(x=math1, y = ..density..)) +
  geom_histogram(binwidth = 20, fill = "grey", color = "black") +
  geom_line(data = df, aes(x = x, y = y), color = "red") +
  labs(x="math score",y="",title = "histogram of math socre in 1st grade")
```



```
qqnorm(data_remove_na$math1, pch = 1, frame = FALSE)
qqline(data_remove_na$math1, col = "steelblue", lwd = 2)
```

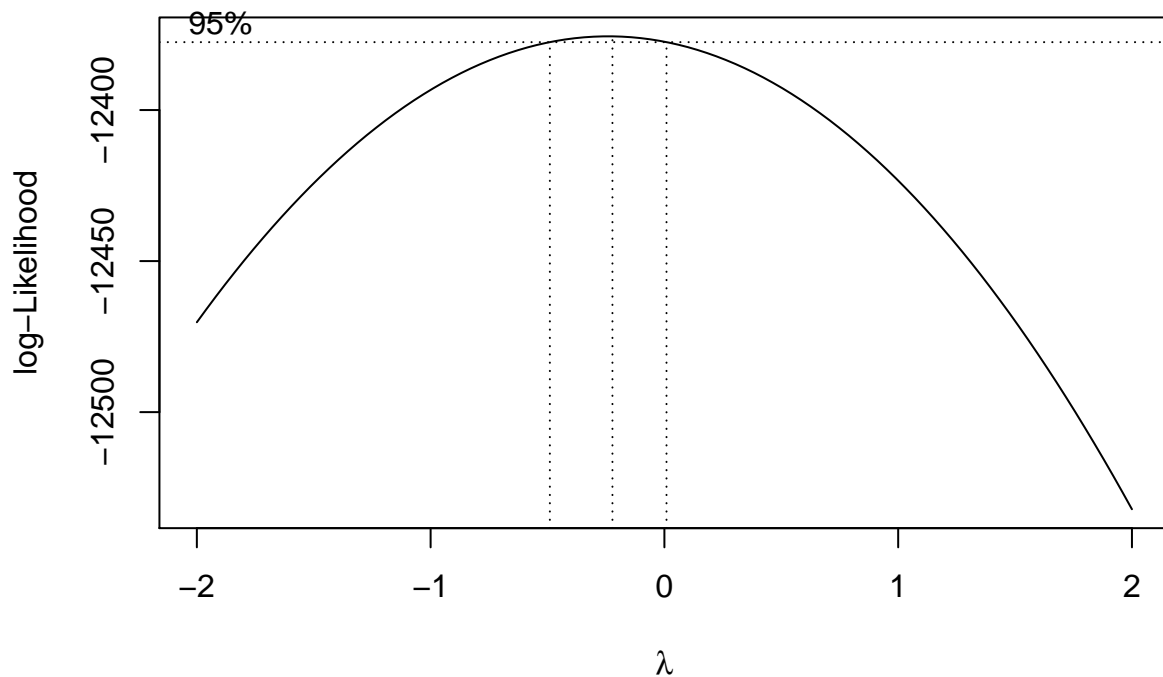


The histogram shows that it seems normal distribution.

The Q-Q plot shows the the distribution of math score is right-skewed.

So we use Box-Cox method to make a transformation on `math1`.

```
library(MASS)
boxcox(math1 ~ star1 , data = data_remove_na)
```



It indicates that we need make a log-transformation for math1.

```
summary(log(data_remove_na$math1))
```

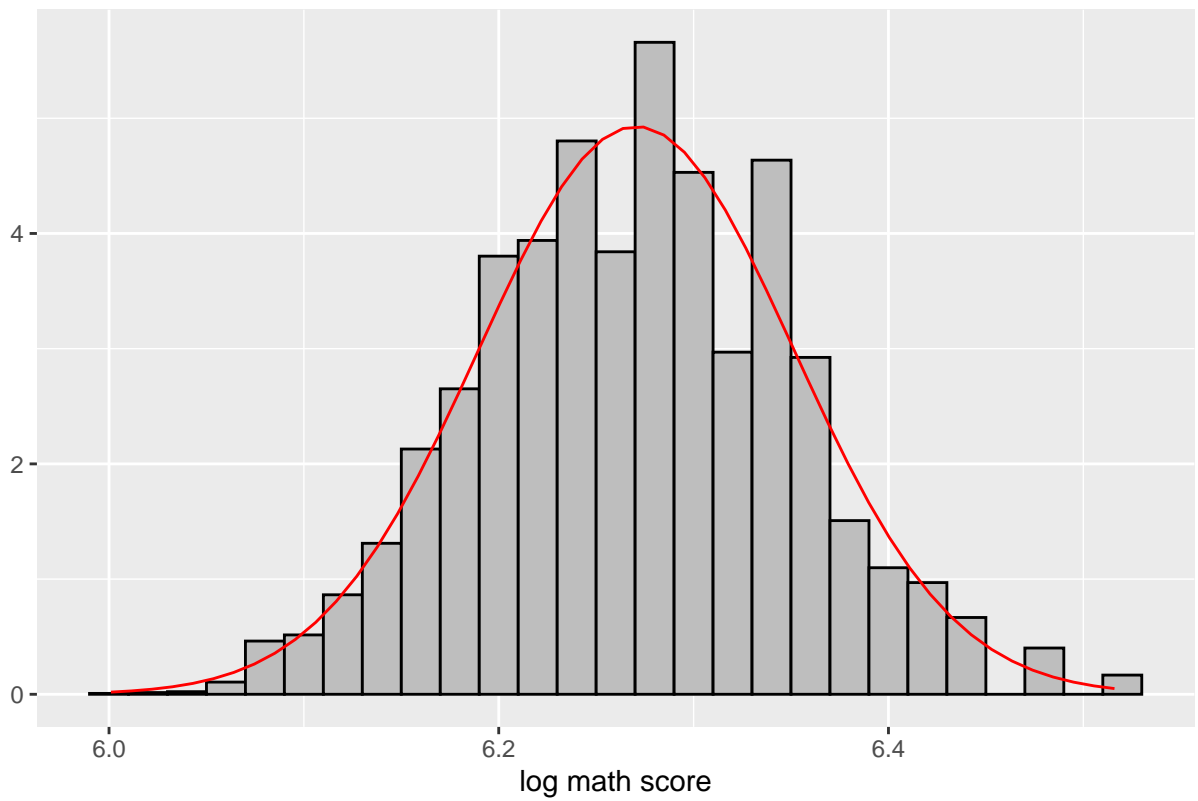
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      6.001  6.215   6.271   6.271  6.323   6.516
```

```
x <- seq(6.001, 6.516, length.out=50)
```

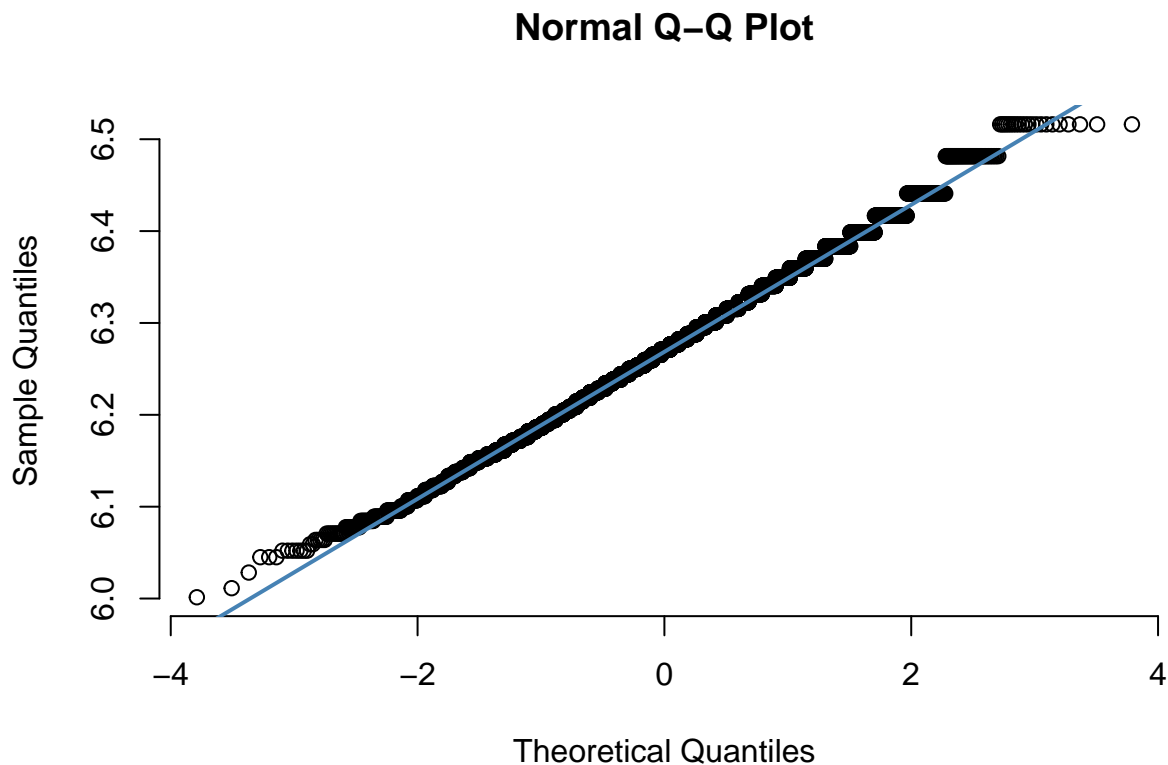
```
df <- with(data_remove_na, data.frame(x = x, y = dnorm(x, mean(log(math1)), sd(log(math1)))))
```

```
ggplot(data_remove_na, aes(x=log(math1), y = ..density..)) +
  geom_histogram(binwidth = 0.02, fill = "grey", color = "black") +
  geom_line(data = df, aes(x = x, y = y), color = "red") +
  labs(x="log math score",y="",title = "histogram of log math socre in 1st grade")
```

histogram of log math socre in 1st grade



```
qqnorm(log(data_remove_na$math1), pch = 1, frame = FALSE)
qqline(log(data_remove_na$math1), col = "steelblue", lwd = 2)
```



The graph shows the distribution of log math score in 1st grade is Normal-like.

Then we calculate the variance of math grade in 1st grade of each class type.

```
library(tidyverse)

data_remove_na %>%
  group_by(star1) %>%
  summarize(var_math1 = var(log(math1), na.rm = T))
```

```
## # A tibble: 3 x 2
##   star1      var_math1
##   <fct>      <dbl>
## 1 regular    0.00625
## 2 small     0.00666
## 3 regular+aide 0.00650
```

The result shows that they are very small and nearly equal to each other. Therefore, it is appropriate to build our model on this dataset.

### Step5 Fit Model

```
anova.fit<- aov(log(math1)~star1,data=data_remove_na)
summary(anova.fit)
```

```
##           Df Sum Sq Mean Sq F value Pr(>F)
## star1      2   0.68   0.3391   52.56 <2e-16 ***
```



```
## Residuals    6597   42.55   0.0065
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

anova.fit$coefficients
```

```
##      (Intercept)      star1small star1regular+aide
##      6.26079457      0.02499081      0.00812069
```

From the result, the fitted model we get is:

$$\log \hat{Y}_{ij} = 6.2608 + 0.0250X_{2,ij} + 0.0081X_{3,ij}$$

with means when the type is regular, the estimate math score is  $e^{6.2608} = 523.6377$ ; when the type is small, the estimate math score is  $e^{6.2608+0.0250} = 536.8936$ ; when the type is regular-with-aide, the estimate math score is  $e^{6.2608+0.0081} = 527.8964$ .

The following is a ANOVA table for this model.

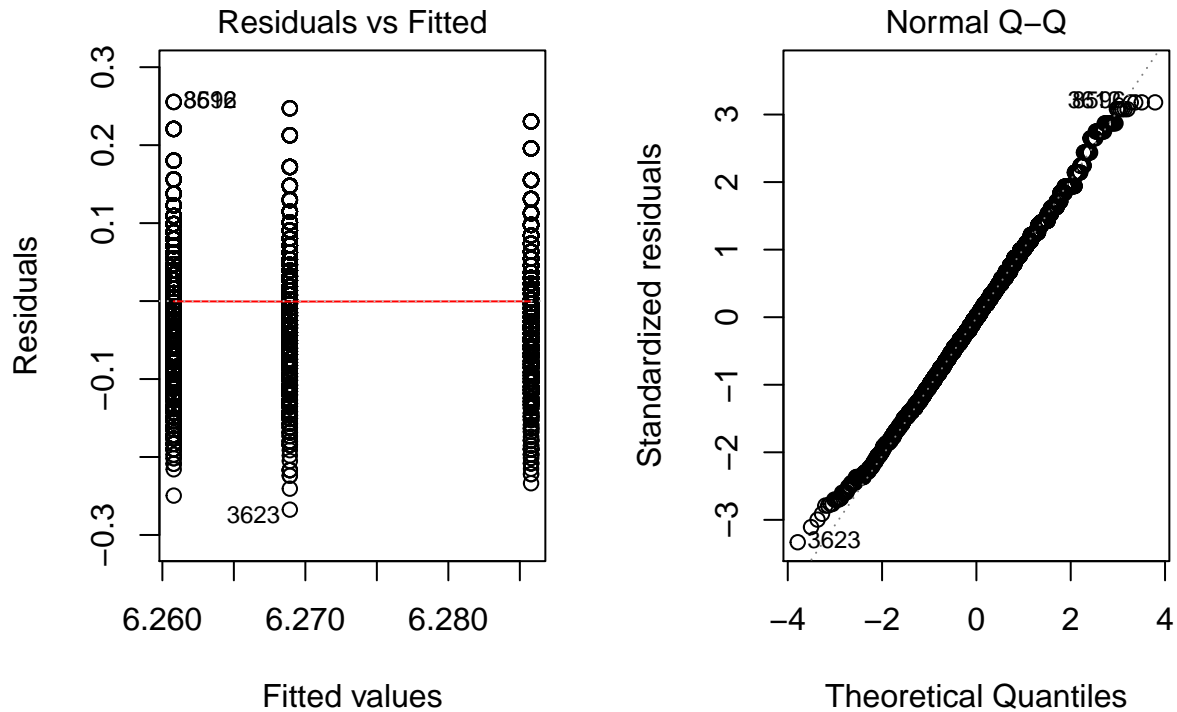
Source of Variation	Sum of Squares	Degrees of Freedom	MS
Between treatments	SSTR = 0.68	2	MSTR = 0.3391
Within treatments	SSE = 42.55	6597	MSE = 0.0065
Total	SSTO = 43.23	6599	

## Step6 model diagnostic and sensitivity analysis

Recalling the above assumptions, there are three things we need to check: normality, equal variance and independence.

By Q-Q plot, we can check normality. And By residuals vs fitted value plot, we can check equal variance.

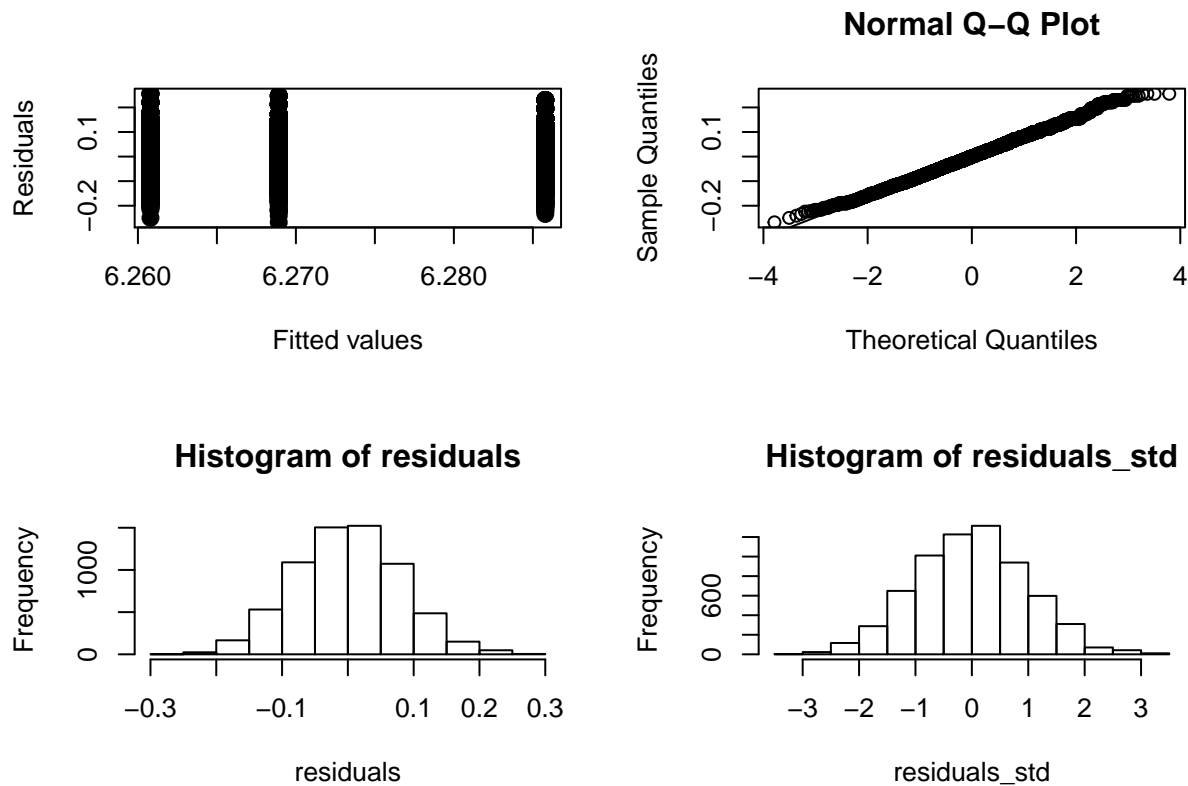
```
par(mfrow=c(1,2))
plot(anova.fit, which = c(1,2))
```



The first figure indicates equal variance, and the second nearly linear figure supports normality.

shijianversion

```
par(mfrow=c(2,2))
residuals <- anova.fit$residuals
##Plot the residuals (or the other two versions) against fitted values
plot(anova.fit$fitted.values, anova.fit$residuals,
     type = "p", pch=16, cex=1.5, xlab="Fitted values", ylab="Residuals")
#QQplot
qqnorm(residuals); qqline(residuals)
#residuals
hist(residuals)
#studentized residuals
residuals_std <- rstudent(anova.fit)
hist(residuals_std)
```



From the scatterplot of residuals vs fitted values, the residuals are divided into three groups and among each group, these residuals are around the zero, which means that the average residuals equals to zero. According to the histogram of residuals and studentized residuals, we can find that the distribution of the residuals of the fitted model approximates to the normal distribution. Besides, the same conclusion can be obtained by checking the Q-Q Plot of the residuals. Therefore, we can confirm that the residuals of the model are normally distributed.

We now turn to formal tests of the equality of variances. First, we calculate the variances for each type of class and find that the variances of three types of class are close to each other.

```
# Calculate the variances for each group:
(vars = tapply(data_remove_na$math1, data_remove_na$star1, var))
```

```
##      regular      small regular+aide
## 1734.825    1945.082    1837.493
```

Then, because the sample sizes of the three types of class are not same, we choose two formal tests, which are Bartlett test and Levene test, to check the equality of model variances.

```
data_remove_na$residuals <- residuals
#bartlett test
bartlett.test(residuals ~ star1, data = data_remove_na)
```

```
##
## Bartlett test of homogeneity of variances
##
## data:  residuals by star1
## Bartlett's K-squared = 2.3004, df = 2, p-value = 0.3166
```

```
#levene test
leveneTest(residuals ~ star1, data = data_remove_na)
```

```
## Levene's Test for Homogeneity of Variance (center = median)
##           Df F value Pr(>F)
## group      2  1.1836 0.3062
##           6597
```

From the two tests, both of the P-values are much larger than 0.05, which means that we can not reject the null hypothesis: the variances of the model are equal.

In conclusion, we confirm that our model satisfies the normality assumption.

**Sensitivity Analysis** In order to test the sensitivity of our model, we decide to relax the assumption of our model. To be specific, we want to figure out that whether the influence of class size still exists even if the data is not normally distributed. Thus, we conduct the nonparametric tests as follows, which are The rank test and Kruskal-Wallis test.

```
#rank test
data_remove_na$rank <- rank(data_remove_na$math1)
summary(aov(rank ~ star1, data = data_remove_na))
```

```
##           Df      Sum Sq   Mean Sq F value Pr(>F)
## star1      2 3.556e+08 177779655   49.72 <2e-16 ***
## Residuals 6597 2.359e+10   3575611
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
#kruskal test
kruskal.test(math1 ~ star1, data = data_remove_na)
```

```
##
## Kruskal-Wallis rank sum test
##
## data:  math1 by star1
## Kruskal-Wallis chi-squared = 97.993, df = 2, p-value < 2.2e-16
```

The results both show that the math scores of the different types of class are different at 99% confident level. So, even if the data was not normally distributed, there would still be influence of class size. In a word, our one-way anova model is reasonable in this case.

#step6 hypothese test

```
TukeyHSD(anova.fit)
```

```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = log(math1) ~ star1, data = data_remove_na)
##
## $star1
##           diff          lwr          upr          p adj
## small-regular 0.02499081 0.019236297 0.03074533 0.0000000
## regular+aide-regular 0.00812069 0.002637091 0.01360429 0.0015104
## regular+aide-small -0.01687012 -0.022778294 -0.01096196 0.0000000
pairwise.t.test(log(data_remove_na$math1),data_remove_na$star1,p.adj = "bonf")
```

```
##
## Pairwise comparisons using t tests with pooled SD
##
## data: log(data_remove_na$math1) and data_remove_na$star1
##
##          regular small
## small      < 2e-16 -
## regular+aide 0.0016  7.1e-11
##
## P value adjustment method: bonferroni
```

```
library(agricolae)
scheffe.test(anova.fit,"star1", group=TRUE,console=TRUE)
```

```
##
## Study: anova.fit ~ "star1"
##
## Scheffe Test for log(math1)
##
## Mean Square Error : 0.006450156
##
## star1, means
##
##          log.math1.      std      r      Min      Max
## regular          6.260795 0.07904431 2507 6.011267 6.516193
## regular+aide      6.268915 0.08062925 2225 6.001415 6.516193
## small            6.285785 0.08161399 1868 6.052089 6.516193
##
## Alpha: 0.05 ; DF Error: 6597
## Critical Value of F: 2.997093
##
## Groups according to probability of means differences and alpha level( 0.05 )
##
## Means with the same letter are not significantly different.
##
##          log(math1) groups
## small          6.285785      a
## regular+aide    6.268915      b
## regular          6.260795      c
```

For task 7, we choose three methods to test the difference in the math scaled score in 1st grade across students in different class types. We use Tukey's Procedure, Bonferroni's Procedure and Scheffe's procedure. For Tukey's Procedure, all the p values are less than 0.05, there is statistically significant among three factors. For Bonferroni's Procedure, we get the same result as Tukey's Procedure. However, for Scheffe's procedure, we get the different result. It shows that Means with the same letter are not significantly different.