

Jack Dumanski, Ryan O'Neill

Dr. Perry

STAT 351 001

7 November 2025

Exploratory Data Analysis

Our data set (MathScores.csv) is surrounding students and their success rate in their math classes based on several factors on the personal and academic lives of the students. The target variable is categorical and called "G4" which has entries "F" – Fail, "NI" – Needs Improvement, "S" – Satisfactory, and "E" – Excellent. This variable is related to the student's success in their math class. This model will aim to explain the students' math class success based primarily on variables outside of school performance. To outline each factor, we made plots and calculations below for all types of variables:

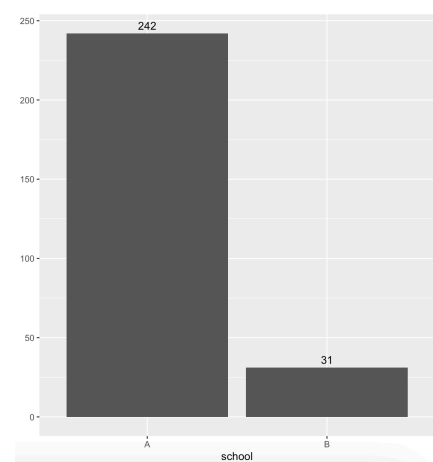
Categorical Variables:

```
library(ggplot2)
```

School:

```
ggplot(MathScores, aes(x=school)) + geom_bar() +  
geom_text(stat="count", aes(label=..count..), vjust=-.5)
```

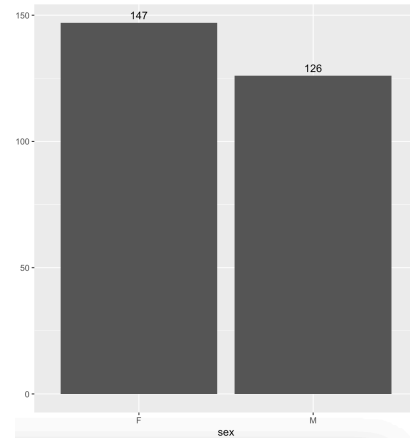
The vast majority of students in the dataset attended school A as opposed to school B.



Sex:

```
ggplot(MathScores, aes(x=sex)) + geom_bar() +  
geom_text(stat="count", aes(label=..count..), vjust=-.5)
```

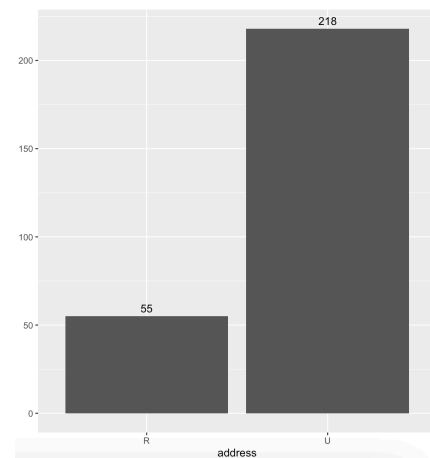
The distribution of Male and Female students is fairly even. There are slightly more female than male students.



Address:

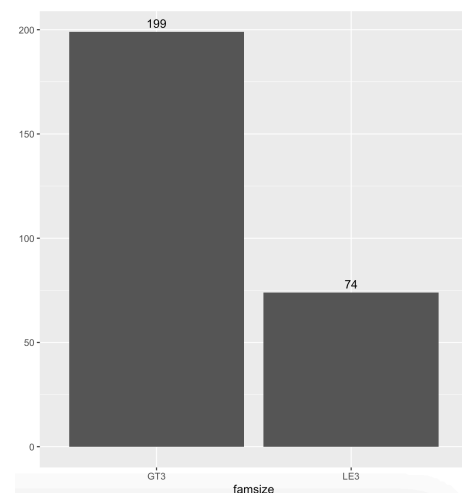
```
ggplot(MathScores, aes(x=address)) + geom_bar() +  
geom_text(stat="count", aes(label=..count..), vjust=-.5)
```

The majority of students in this study are from urban addresses rather than rural.



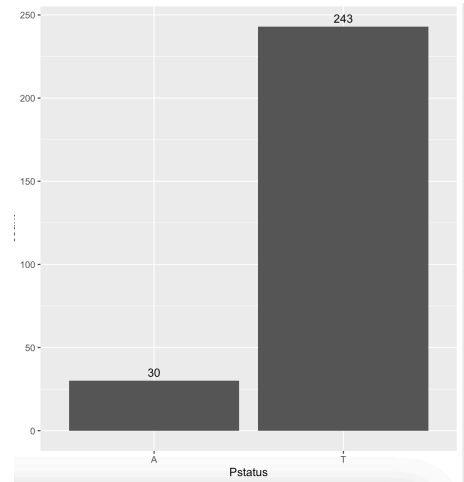
```
Family size: ggplot(MathScores, aes(x=famsize)) +  
geom_bar() + geom_text(stat="count",  
aes(label=..count..), vjust=-.5)
```

The family size for most students is greater than 3, which is around three times the amount of families with size less than 3.



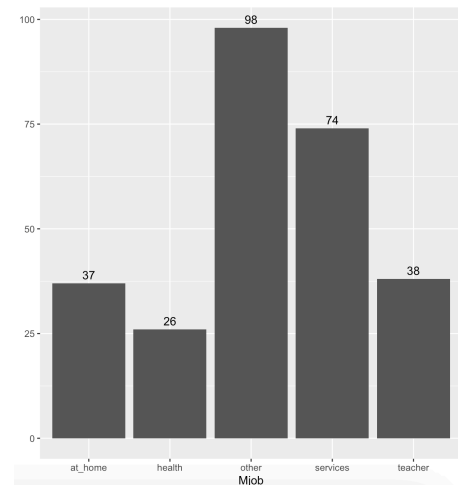
Parent Status: `ggplot(MathScores, aes(x=Pstatus)) +
geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

An overwhelming majority of students in this sample have their parents together rather than apart.



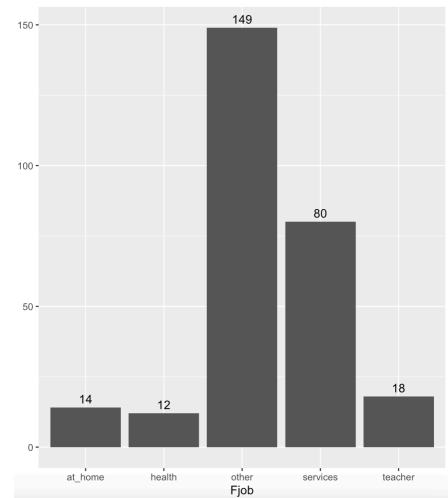
Mother's Job: `ggplot(MathScores, aes(x=Mjob)) +
geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

Out of the common jobs listed (at-home work, healthcare, services, teaching) for the students' mothers, service jobs are the most popular. However, a large number of students' mothers do not have a common job.



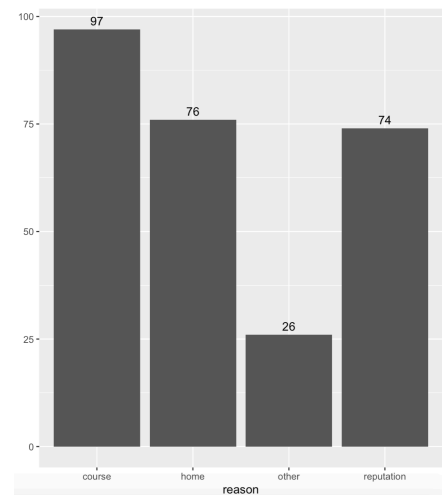
Father's job: `ggplot(MathScores, aes(x=Fjob)) +
geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

Students' fathers tend to have the most jobs in services rather than the other listed common jobs, but the majority of jobs fall in the "other" category.



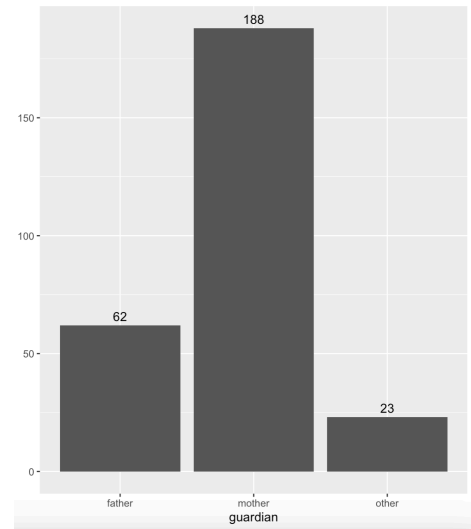
Reason for choosing school: `ggplot(MathScores,
aes(x=reason)) + geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

The most common reason students chose their school is due to their course offerings, with proximity to home and the school's reputation not too far behind.



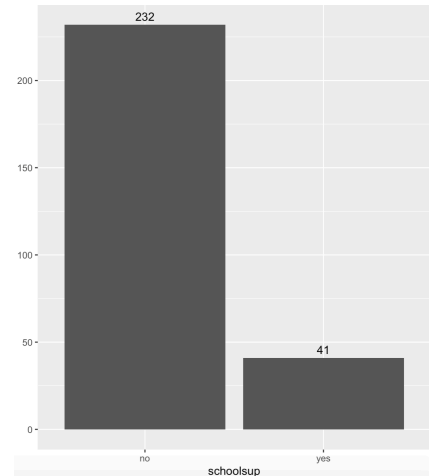
Guardian choice: `ggplot(MathScores, aes(x=guardian)) +
geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

The majority of students list their mother as their primary guardian, around 3 times as common as they list their father.

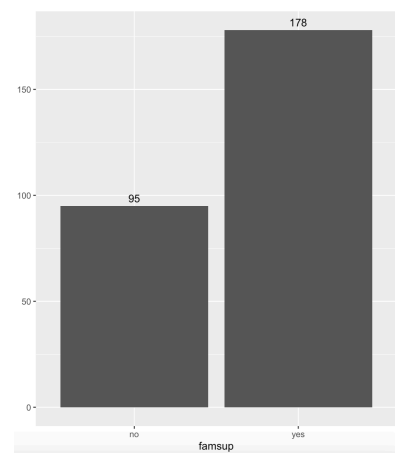


School support: `ggplot(MathScores, aes(x=schoolsup)) +
geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

The majority of students do not have additional school support, but 41 students do.



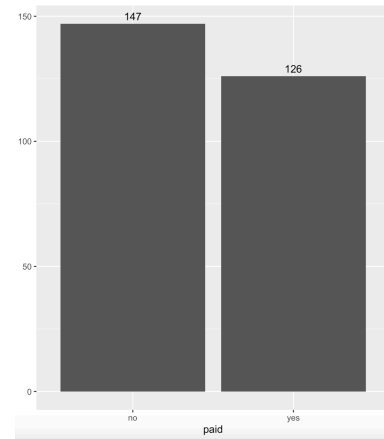
Family support: `ggplot(MathScores, aes(x=famsup)) +
geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`



The majority of students do have family educational support, and it is around twice as often that students do have family educational support than those who don't.

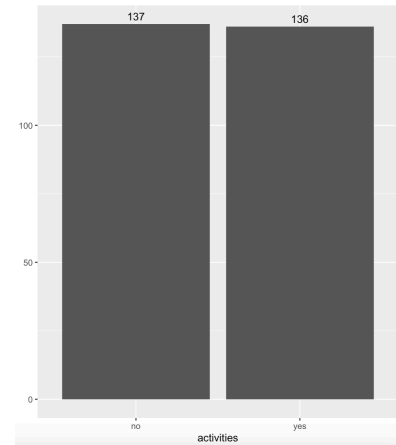
Extra math classes: `ggplot(MathScores, aes(x=paid)) +
geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

There is a roughly even split between students who have extra paid Math classes and those who don't.



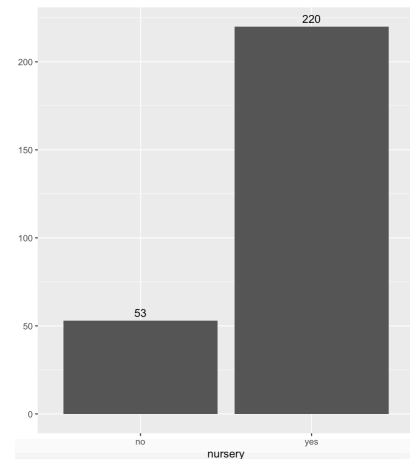
Activities: `ggplot(MathScores, aes(x=activities)) +
geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

There is a nearly even split of students who participate in extra-curricular activities and those who don't.



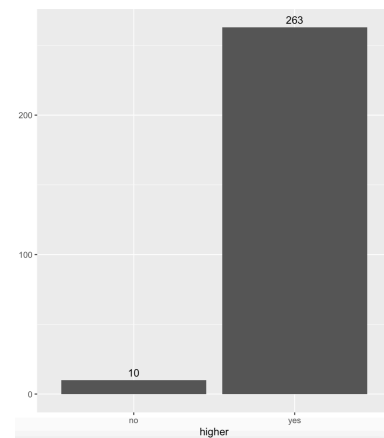
Nursery school: `ggplot(MathScores, aes(x=nursery)) +
geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

The majority of students attended nursery school, but there is still a strong population of students who didn't.



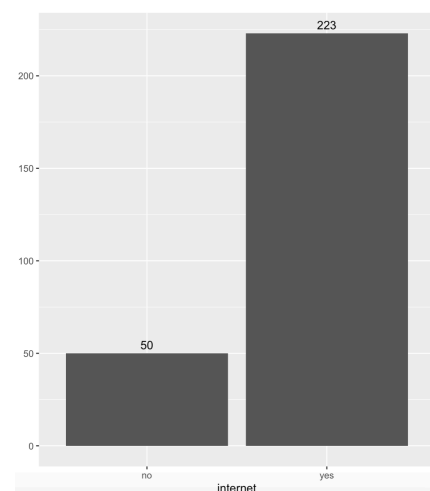
Pursuing higher education: `ggplot(MathScores, aes(x=higher)) + geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

Nearly every student in this study has goals to pursue higher education.



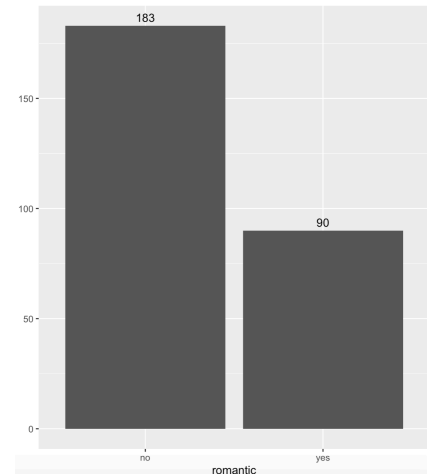
Internet access: `ggplot(MathScores, aes(x=internet)) +
geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

A large majority of students in this study have internet access at home, but there is a strong population of students without home internet access.



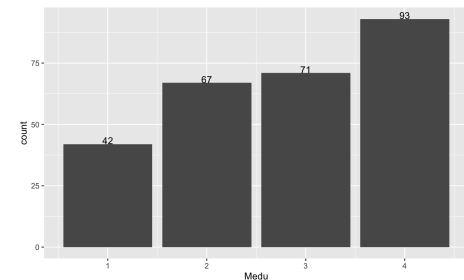
Romantic relationship: `ggplot(MathScores, aes(x=romantic)) + geom_bar() + geom_text(stat="count", aes(label=..count..), vjust=-.5)`

The majority of students in this study do not consider themselves in a romantic relationship, which is around twice as often as those who are in a romantic relationship

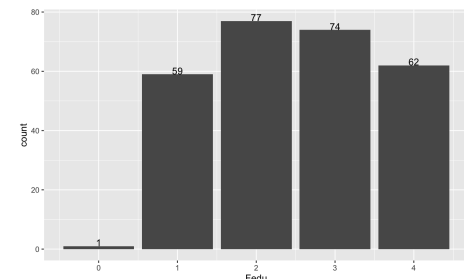


Mother's education: `ggplot(MathScores, aes(x=Medu)) + geom_bar() + geom_text(stat="count", aes(label=..count..), vjust=-.5)`

There is a slight negative skew in the level of mother's education where the highest bar is higher education, then high school, then middle school, then some elementary school.



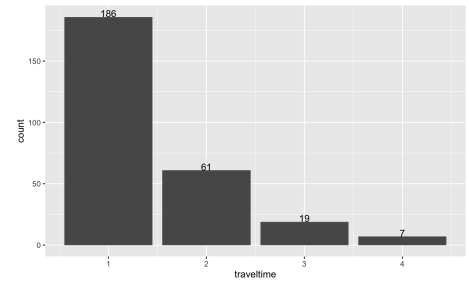
Father's education: `ggplot(MathScores, aes(x=Fedu)) + geom_bar() + geom_text(stat="count", aes(label=..count..), vjust=-.5)`



The bars for father's education are all roughly equal with the same strata as the mother's education, but there is one student whose father had no education at all.

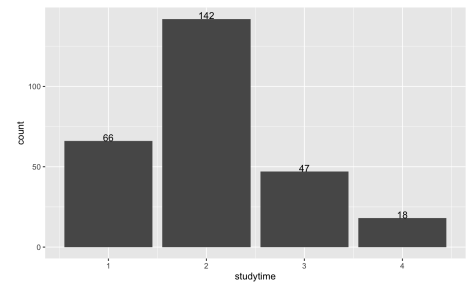
Travelttime: `ggplot(MathScores, aes(x=travelttime)) +
geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

The majority of students are less than 15 minutes away from their school, and there is a positive skew in the data in terms of distance away per student.



Studytime: `ggplot(MathScores, aes(x=studytime)) +
geom_bar() + geom_text(stat="count",
aes(label=..count..), vjust=-.5)`

The most common amount of time spent studying per student is 2-5 hours per week, with 0-2 and 5-10 being the next highest bars, and the fewest students studying 10+ hours.



Numerical Variables:

Failed classes: `ggplot(aes(x = failures))`

+

`geom_bar()`

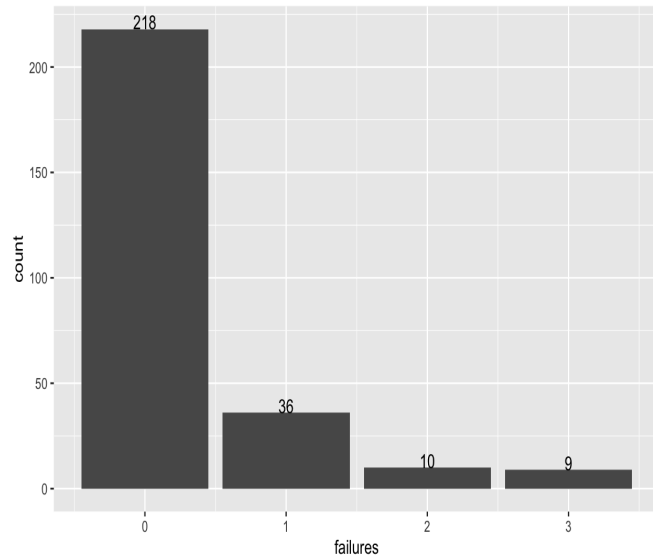
The majority of students did not fail any classes, and this distribution has a strong positive skew, with very few students failing any classes, especially more than 1.

`summary(scores$failures)`

`mean(scores$failures)`

`sd(scores$failures)`

The average amount of classes failed by students is 0.3 with a standard deviation of 0.69.



Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.000	0.000	0.000	0.304	0.000	3.000

```
Absences: ggplot(aes(x = absences)) +  
  geom_histogram(bins = 10)
```

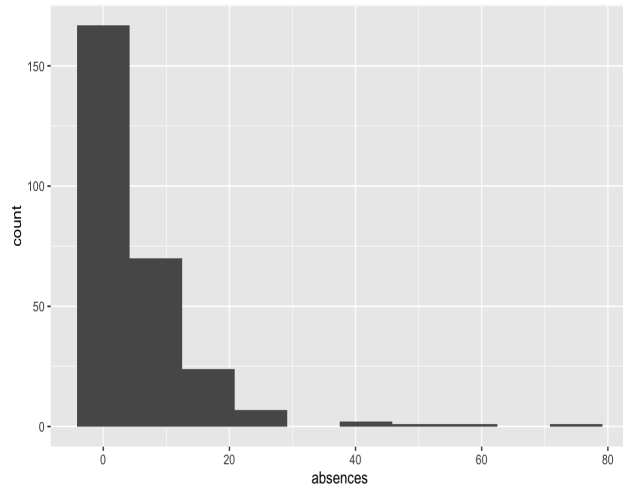
This distribution of absences per student has a strong positive skew, and students usually missed no more than 10 classes, and virtually no students missed more than 20.

```
summary(scores$absences)
```

```
mean(scores$absences)
```

```
sd(scores$absences)
```

The average absences per student is 6.01 with a standard deviation of 8.74, which is influenced by the large outliers.



Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.000	0.000	4.000	6.011	8.000	75.000

6.010989

8.739483

```
Age: ggplot(aes(x = age)) + geom_bar()
+
  geom_text(stat = "count", aes(label =
    ..count..), vjust = 0 )
```

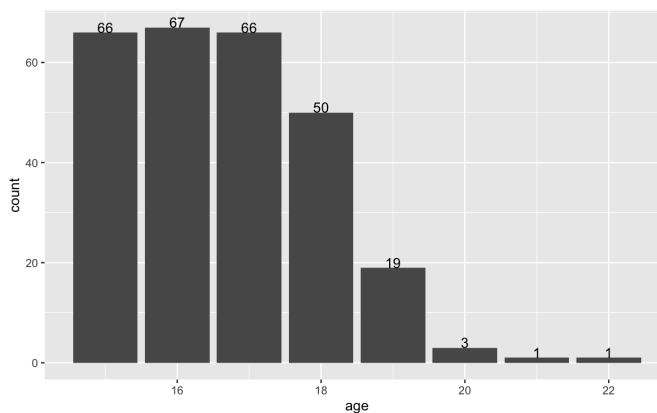
The majority of the students polled were from ages 15-18, which is traditional for high school students. A few outliers were older, up to 22 years old.

```
summary(scores$age)
```

```
mean(scores$age)
```

```
sd(scores$age)
```

The average age of students polled is 16.65 with a standard deviation of 1.35, which means 68% of students are aged between 15.3 and 18.



Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
15.00	16.00	17.00	16.66	18.00	22.00

```
16.65934
```

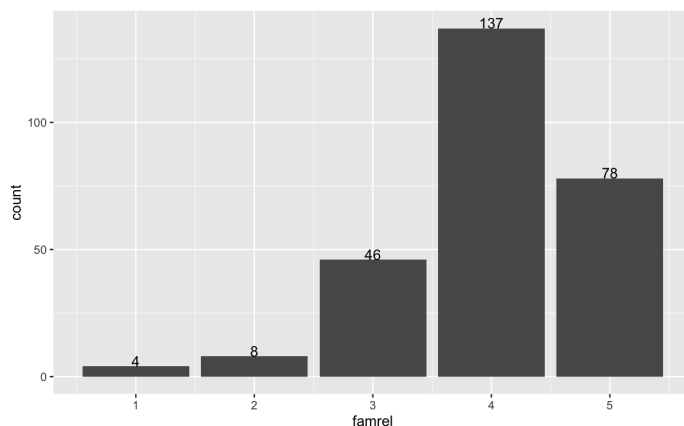
```
1.346726
```

```
Family Relationship: ggplot(aes(x =
Famrel)) + geom_bar() +
  geom_text(stat = "count", aes(label =
  ..count..), vjust = 0 )
```

The majority of students had a 3+ relationship with their family on a 1-5 scale of quality of their family relationship. This distribution is negatively skewed.

```
summary(scores$famrel)
mean(scores$famrel)
sd(scores$famrel)
```

The average rating of students' family relationship is 4.01/5, with a 0.84 standard deviation.



Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.000	4.000	4.000	4.015	5.000	5.000
4.014652					
0.8400398					

```
Free time: ggplot(aes(x = freetime)) +
geom_bar() +
  geom_text(stat = "count", aes(label =
    ..count..), vjust = 0 )
```

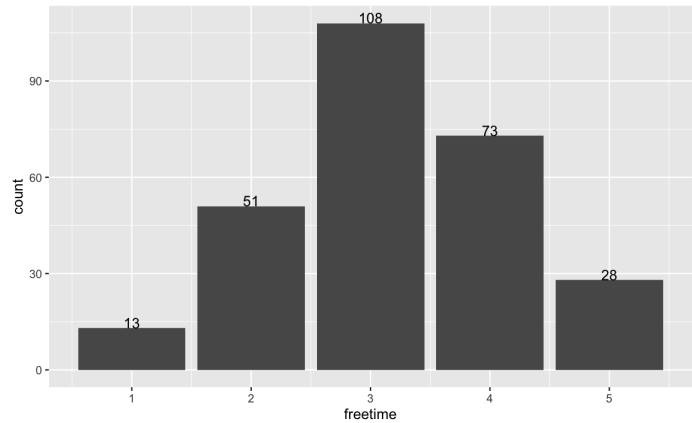
This graph is roughly normally distributed, with the majority of students centered around 3 amount of freetime.

```
summary(scores$freetime)
```

```
mean(scores$freetime)
```

```
sd(scores$freetime)
```

The average rating of students' free time from 1-5 is 3.2 with a standard deviation of 1.01.



Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.00	3.00	3.00	3.19	4.00	5.00

```
Going out: ggplot(aes(x = goout)) +
  geom_bar() +
  geom_text(stat = "count", aes(label =
    ..count..), vjust = 0 )
```

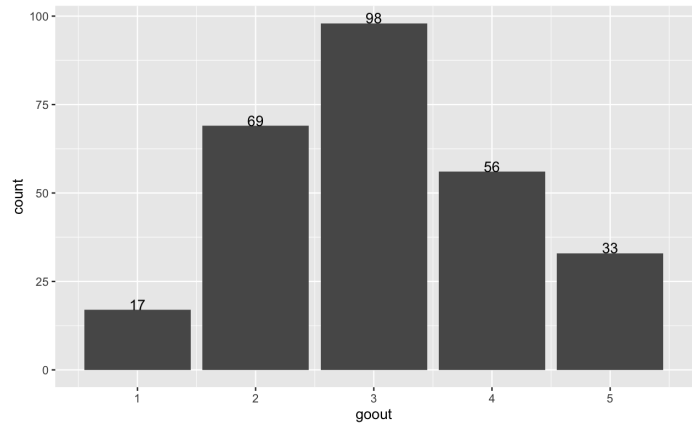
The distribution for students' going out time rating is approximately normal centered around 3.

```
summary(scores$goout)
```

```
mean(scores$goout)
```

```
sd(scores$goout)
```

The average rating out of 5 for students' going out time is 3.07 with a standard deviation of 1.09.



Min. 1st Qu. Median Mean 3rd Qu. Max.

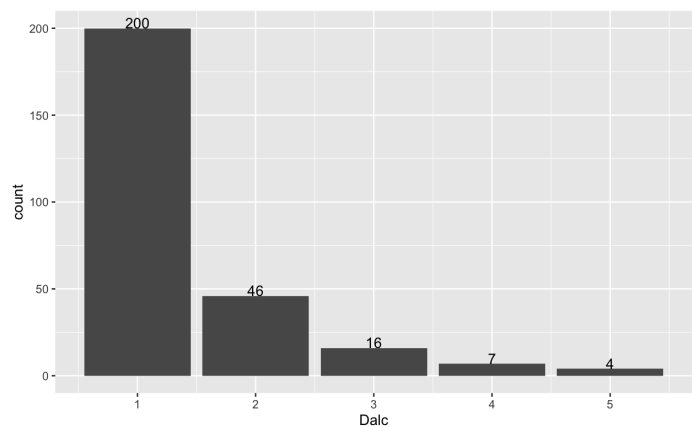
1.00 2.00 3.00 3.07 4.00 5.00

3.069597

1.090867

Work day alcohol consumption:

```
ggplot(aes(x = Dalc)) + geom_bar() +
  geom_text(stat = "count", aes(label =
    ..count..), vjust = 0 )
```



This distribution of rating of workday alcohol consumption is strongly positively skewed.

```
summary(scores$Dalc)
```

```
mean(scores$Dalc)
```

```
sd(scores$Dalc)
```

The average rating out of 5 of alcohol consumption on a workday is 1.42 with a standard deviation of 0.832.

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.000	1.000	1.000	1.421	2.000	5.000

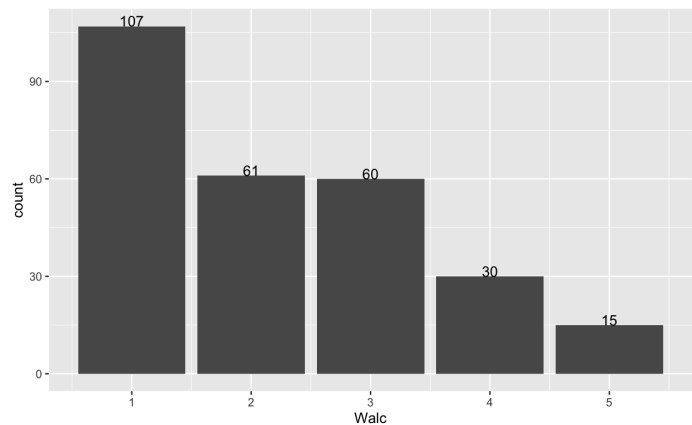
Weekend Alcohol consumption:

```
ggplot(aes(x = Walc)) + geom_bar() +  
  geom_text(stat = "count", aes(label =  
    ..count..), vjust = 0 )
```

This distribution has a strong positive skew, with the majority of students rating 3 or less for their weekend alcohol consumption.

```
summary(scores$Walc)
```

```
mean(scores$Walc)
```



Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.000	1.000	2.000	2.212	3.000	5.000


```
sd(scores$Walc)
```

The average weekend alcohol consumption rating out of 5 is 2.21 with a standard deviation of 1.227.

```
Health rating: ggplot(aes(x = health)) +  
geom_bar() +  
  geom_text(stat = "count", aes(label =  
..count..), vjust = 0 )
```

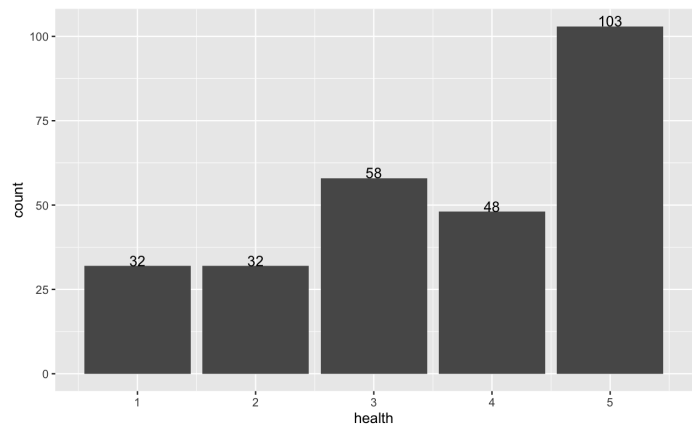
This distribution is largely negatively skewed, with the majority of students rating their health to be “very high” or close to it.

```
summary(scores$health)
```

```
mean(scores$health)
```

```
sd(scores$health)
```

The average health rating for students is 3.58 with a standard deviation of 1.39, showing high average health.



Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.000	3.000	4.000	3.579	5.000	5.000

Conclusion:

Most of the variables in this dataset are categorical or based on perceived ratings on 1-5 scales, so a tree diagram may be the best method to form a model for this data. We may first form a tree with every variable, factored, in a tree model but prune the tree until its predictive ability is adequate enough, or we may begin with variables we think are most important and then test for significance in each additional variable. The most pertinent variables we noticed are Studytime, Schoolsup, Famsup, Paid, Internet, Famrel, Absences. To us, these represented direct connections to success in school or psychological effects that may prohibit students from succeeding. Again, since the data has very few numerical variables, the tree will do well to split the data based on categorical strata, for example $\text{health} \leq 3$. This is nonparametric because we can't make assumptions about the underlying distributions of these variables and simply add them as categories into our tree.