

## Exercise 9 - Students Surevy Example

Ninad Patkhedkar

2020-10-01

As a data science intern with newly learned knowledge in skills in statistical correlation and R programming, you will analyze the results of a survey recently given to college students. You learn that the research question being investigated is: “Is there a significant relationship between the amount of time spent reading and the time spent watching television?” You are also interested if there are other significant relationships that can be discovered? The survey data is located in this `StudentSurvey.csv` file.

```
##   TimeReading TimeTV Happiness Gender
## 1           1     90    86.20      1
## 2           2     95    88.70      0
## 3           2     85    70.17      0
## 4           2     80    61.31      1
## 5           3     75    89.52      1
## 6           4     70    60.50      1
```

a) Use R to calculate the covariance of the Survey variables and provide an explanation of why you would use this calculation and what the results indicate.

```
cov(student_survey_df)
```

```
##           TimeReading      TimeTV  Happiness      Gender
## TimeReading  3.05454545 -20.36363636 -10.350091 -0.08181818
## TimeTV      -20.36363636 174.09090909 114.377273  0.04545455
## Happiness   -10.35009091 114.37727273 185.451422  1.11663636
## Gender      -0.08181818  0.04545455  1.116636  0.27272727
```

```
gender0_df <- student_survey_df[ which( student_survey_df$Gender == "0"), ]
paste('Covariance for Gender0')
```

```
## [1] "Covariance for Gender0"
```

```
cov(gender0_df)
```

```
##           TimeReading      TimeTV  Happiness      Gender
## TimeReading  3.200 -30.250  -21.4890      0
## TimeTV      -30.250 305.000  248.9050      0
## Happiness   -21.489 248.905  266.9072      0
## Gender       0.000  0.000   0.0000      0
```

```
paste('Covariance for Gender1')
```

```
## [1] "Covariance for Gender1"
```

```
gender1_df <- student_survey_df[ which( student_survey_df$Gender == "1"), ]
cov(gender1_df)
```

```
##           TimeReading    TimeTV Happiness Gender
## TimeReading      3.500 -16.50000  -2.83900      0
## TimeTV          -16.500 104.16667  29.25833      0
## Happiness       -2.839  29.25833 148.23330      0
## Gender           0.000   0.00000   0.00000      0
```

Covariance is a measurement of how closely related two variables are based on a linear relationship. In this example, we received a covariance value of -20 between `TimeTV` and `TimeReading`, meaning that for every hours of reading time a student adds to their daily routine, their daily consumption of television decreases. Their happiness also decreases as covariance value is -10.

Covariance for Gender0 is -30 means Gender0 students likelihood of reduction in TV watching time is more compared to Gender1 students for whom covariance is -16

**b) Examine the Survey data variables. What measurement is being used for the variables? Explain what effect changing the measurement being used for the variables would have on the covariance calculation. Would this be a problem? Explain and provide a better alternative if needed.**

**TimeReading:** Daily time spent on reading by student. It seems to be in hours.

**TimeTV:** Daily time spent on watching TV by student in minutes.

**Happiness:** Some measure of happiness with unknown unit. It seems to be % value as all values range.

**Gender:** Categorical value. The values 0 and 1 represent male/female or female/male.

Changing the `TimeReading` and `TimeTV` attributes so they would both represent time in minutes would reduce the covariance value.

```
cov(student_survey_df$TimeReading*60, student_survey_df$TimeTV)
```

```
## [1] -1221.818
```

**c) Choose the type of correlation test to perform, explain why you chose this test, and make a prediction if the test yields a positive or negative correlation?**

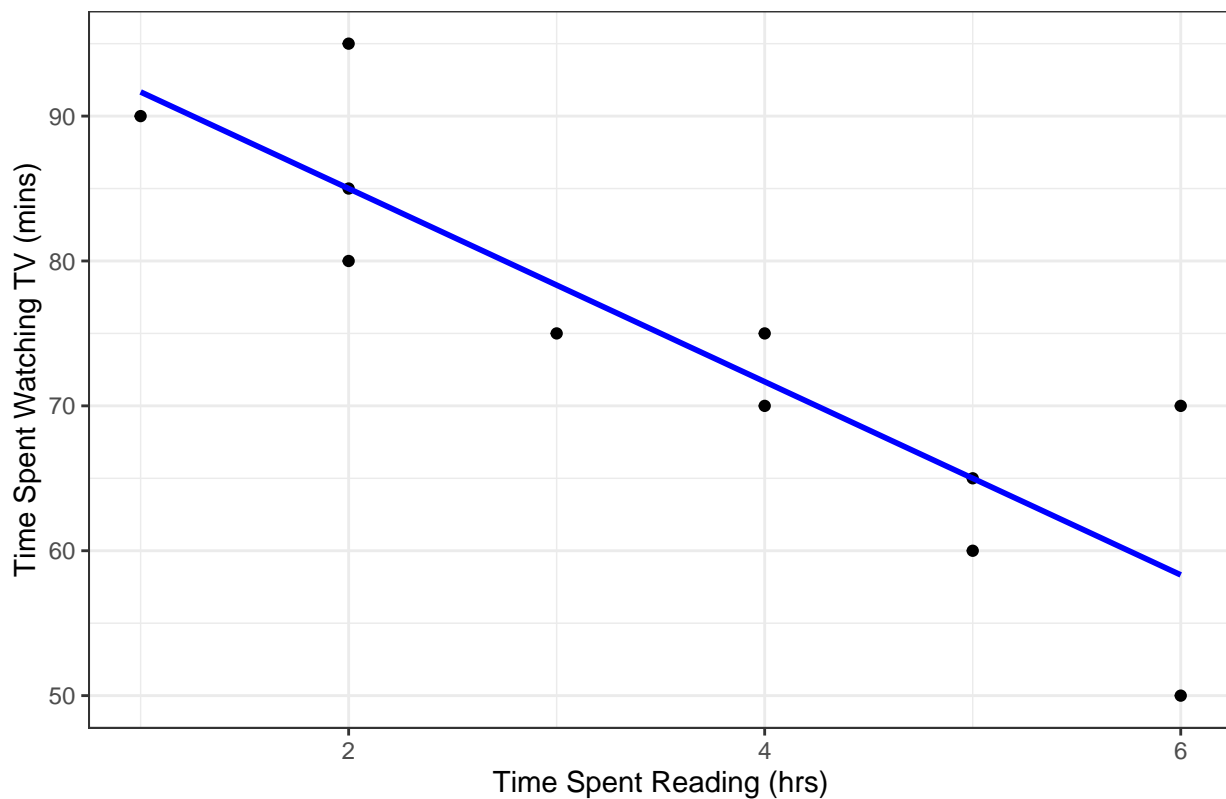
```
##
## Pearson's product-moment correlation
##
## data: student_survey_df$TimeReading and student_survey_df$TimeTV
## t = -5.6457, df = 9, p-value = 0.0003153
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.9694145 -0.6021920
## sample estimates:
##          cor
## -0.8830677
##
## Kendall's rank correlation tau
##
## data: student_survey_df$TimeReading and student_survey_df$TimeTV
## z = -3.2768, p-value = 0.00105
## alternative hypothesis: true tau is not equal to 0
## sample estimates:
```

```
##          tau
## -0.8045404

##
## Spearman's rank correlation rho
##
## data:  student_survey_df$TimeReading and student_survey_df$TimeTV
## S = 419.6, p-value = 0.0001152
## alternative hypothesis: true rho is not equal to 0
## sample estimates:
##          rho
## -0.9072536
```

Pearson, Kendall, and Spearman correlation tests all return p-values less than 5%, indicating a high level of correlation. The correlation values, cor, tau, and rho are also very close to -1, indicating a high negative correlation. We can confirm this visually by plotting the data.

### Student Survey: Reading time and TV time



## d) Correlation analyss

### 1) All variables

```
cor(student_survey_df)
```

```
##          TimeReading      TimeTV  Happiness      Gender
## TimeReading  1.00000000 -0.883067681 -0.4348663 -0.089642146
## TimeTV      -0.88306768  1.000000000  0.6365560  0.006596673
## Happiness   -0.43486633  0.636555986  1.0000000  0.157011838
```

```
## Gender      -0.08964215  0.006596673  0.1570118  1.000000000
```

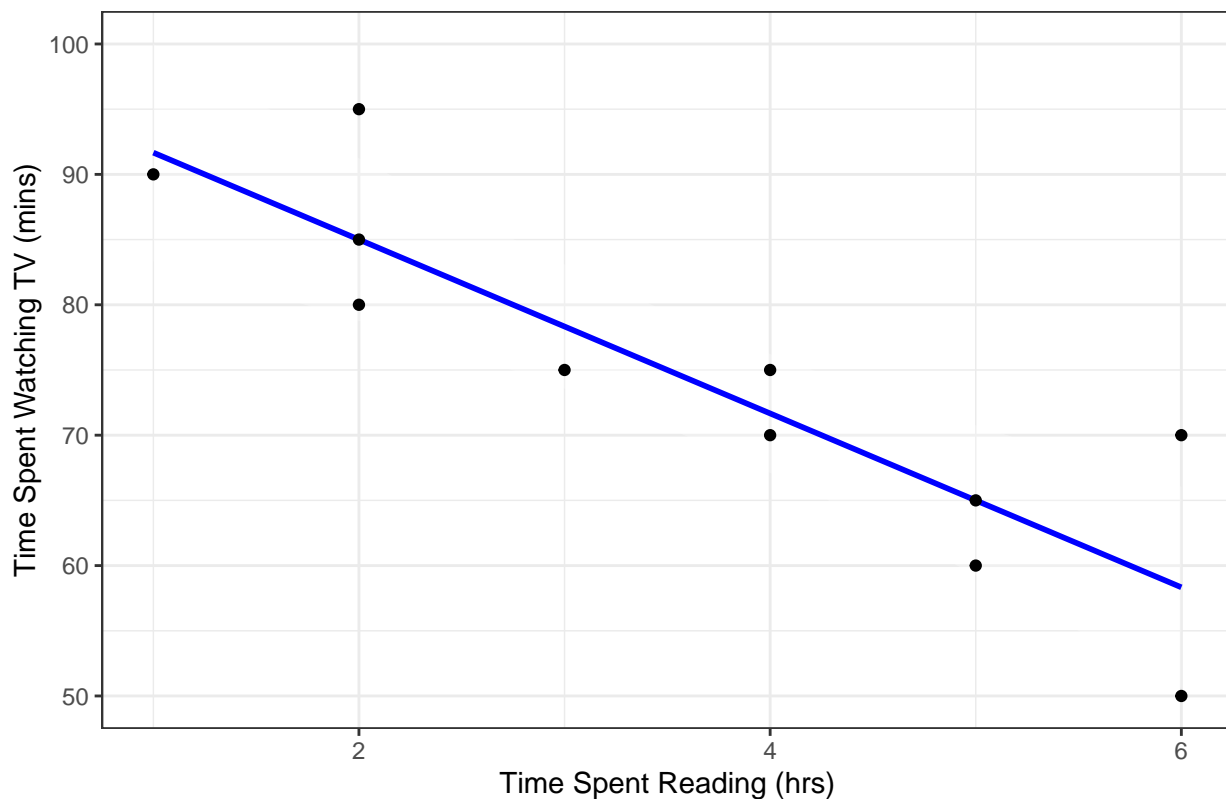
## 2) A single correlation between two a pair of the variables

```
cor.test(formula = ~ student_survey_df$TimeReading + student_survey_df$TimeTV,  
         data = student_survey_df)
```

```
##  
## Pearson's product-moment correlation  
##  
## data: student_survey_df$TimeReading and student_survey_df$TimeTV  
## t = -5.6457, df = 9, p-value = 0.0003153  
## alternative hypothesis: true correlation is not equal to 0  
## 95 percent confidence interval:  
## -0.9694145 -0.6021920  
## sample estimates:  
##      cor  
## -0.8830677
```

```
ggplot(data = student_survey_df,  
       aes(x = TimeReading,  
           y = TimeTV)) +  
geom_smooth(method='lm',  
           formula= y~x,  
           se = TRUE,  
           color = "blue",  
           fill = "white",  
           alpha = 0.3) +  
geom_point() +  
ylab("Time Spent Watching TV (mins)") +  
xlab("Time Spent Reading (hrs)") +  
ggtitle("Student Survey: Reading time and TV time") +  
theme_bw()
```

### Student Survey: Reading time and TV time



3) Repeat your correlation test in step 2 but set the confidence interval at 99%

```
cor.test(formula = ~ student_survey_df$TimeReading + student_survey_df$TimeTV,
         data = student_survey_df, conf.level = 0.99)
```

```
##
## Pearson's product-moment correlation
##
## data: student_survey_df$TimeReading and student_survey_df$TimeTV
## t = -5.6457, df = 9, p-value = 0.0003153
## alternative hypothesis: true correlation is not equal to 0
## 99 percent confidence interval:
## -0.9801052 -0.4453124
## sample estimates:
## cor
## -0.8830677
```

4) Describe what the calculations in the correlation matrix suggest about the relationship between the variables. Be specific with your explanation.

For correlation, the closer the value is to 1 or -1, the stronger the variables are correlated. We see a correlation score of -0.88 between **TimeTV** and **TimeReading**, indicating that these two attributes are highly negatively correlated. **Gender** isn't strongly correlated with any other attributes. **Happiness** has a slight positive correlation with **TimeTV** and a slight negative correlation with **TimeReading**. Perhaps they should pick better books?

e) Calculate the correlation coefficient and the coefficient of determination, describe what you conclude about the results.

```
cor(student_survey_df$TimeReading, student_survey_df$TimeTV)
```

```
## [1] -0.8830677
```

```
lm_fit <- lm(TimeReading ~ TimeTV , data = student_survey_df)
summary(lm_fit)$r.squared
```

```
## [1] 0.7798085
```

The correlation coefficient returned is -0.883, indicating a high level of negative correlation. The coefficient of determination, which is the correlation coefficient squared, is 0.779, meaning that 77.9% of our data falls into our expected variance.

f) Based on your analysis can you say that watching more TV caused students to read less? Explain.

Based on the p-values we received from the Pearson, Spearman, and Kendall tests plus the correlation values we calculated, we can say with a high confidence that students who watch more TV spend less time reading.