

Lab 01 - Review of Statistical Concepts

PSC-103B

Marwin Carmo

2025-01-05

Central tendency: Mean, Median, and Mode

There are many measures of central tendency, but these 3 are the most common. They are used to describe a distribution of observations (e.g., all the grades on an exam) in one number that best represents that distribution. Let's see how this works.

First, let's create some variables! Suppose we asked a bunch of UC Davis students how many hours per week they spent watching Netflix, and how many hours they spent exercising during Winter break. We record their answers in two separate vectors:

```
netflix <- c(2, 6, 1, 7, 2, 4, 11, 40, 7, 0, 3, 4, 5, 2, 15)
exercise <- c(2, 2, 6, 2, 12, 45, 8, 3, 2, 6, 4, 0, 1, 3, 0)
```

How many observations are in each variable?

```
length(netflix)
```

```
## [1] 15
```

```
length(exercise)
```

```
## [1] 15
```

We have 15 people in our dataset now. Let's take a look at the average time each student spent on these activities:

```
mean(netflix, na.rm = TRUE) # use the argument na.rm = TRUE to ignore missing values
```

```
## [1] 7.266667
```

```
mean(exercise, na.rm = TRUE)
```

```
## [1] 6.4
```

Unsurprisingly, students exercised less than they watched netflix, on average. But is the mean a good representation of these data? Check out that person who looks like their job is watching Netflix. 40 hours a week? What about that athlete who exercised 45 hours per week over the break?

When we have outliers, sometimes the median is a better representation of the data we have. Remember, the median is the middle value of your data, after you have ordered it.

```
median(exercise, na.rm = TRUE)
```

```
## [1] 3
```

That's pretty different from what we had before!

Sometimes, we can't do arithmetic on the data we have. For example, if we had asked our 15 participants what their favorite flavor of ice cream was, we would not be able to describe that distribution using a mean or a median. That's when the Mode is useful: The mode is just the most frequent value. It's not used very often so R doesn't have a function for that. But you can use another very useful function to find the mode: `table()`

```
table(netflix)
```

```
## netflix
##  0  1  2  3  4  5  6  7 11 15 40
##  1  1  3  1  2  1  1  2  1  1  1
```

`table()` gives you the number of times each element shows up in an object. by looking at the results here, we can see that 2 shows up 3 times, so it's the mode of `netflix`.

To make it easier to see, you can use the `sort` function on the result of the `table` function to order it. Here's that for `exercise`. What's the mode?

```
sort(table(exercise))
```

```
## exercise
##  1  4  8 12 45  0  3  6  2
##  1  1  1  1  1  2  2  2  4
```

Spread: Variance and Standard Deviation

Imagine that I told you the mean number of hours students spent exercising each week over the break was 6.4 hours. You would know something about these students' exercise habits, but not a lot. Do all the students exercise about the same? Or do some students exercise a lot while others don't? These are the types of questions that measures of spread try to tackle: how are the observations spread out around the mean or median?

We're gonna look at two different kinds that are related: variance and standard deviation.

To get the variance:

1. Calculate the mean:

```
mean(exercise)
```

```
## [1] 6.4
```

2. Find the distance from each observation to the mean. R can do this for us pretty easily, because it is used to working with vectors. So when we do this:

```
exercise - mean(exercise)
```

```
## [1] -4.4 -4.4 -0.4 -4.4  5.6 38.6  1.6 -3.4 -4.4 -0.4 -2.4 -6.4 -5.4 -3.4 -6.4
```

It will subtract 6.4 from each observation in `exercise`. Let's save that.

```
diffs <- exercise - mean(exercise)
```

3. Square the differences. Just as we did before, this is also very easy in R:

```
diffs_sq <- diffs^2
```

4. Sum everything and divide by N-1:

```
sum(diffs_sq)/14
```

```
## [1] 124.4
```

The variance of `exercise` is understandably a large number. Remember how that one athlete was very far from the mean? We can check our answers using `var()`

```
var(exercise)
```

```
## [1] 124.4
```

To get the standard deviation, we just get the square root of the variance:

```
ex_var <- var(exercise)
sqrt(ex_var)
```

```
## [1] 11.15347
```

Or, R has a `sd()` function:

```
sd(exercise)
```

```
## [1] 11.15347
```

We often prefer to use the standard deviation, because its units are the same units of our variable, unlike variance, which is units²

Relationships between variables: correlation and covariance

Remember that in our imaginary example, we asked each person about both their exercise and netflix activity over the break. Therefore, each person gave us two bits of information. Let's organize our data into a dataframe to better keep track of it.

```
df <- data.frame(Netflix = netflix,
                  Exercise = exercise,
                  stringsAsFactors = FALSE)
```

Here, the capitalized names are the variable names, and I'm assigning the objects `netflix` and `exercise` to those variables. You can look at `df` by clicking on it, using `View()`, typing it in the console, or using functions like `head()` and `tail()`.

```
df
```

```
##      Netflix Exercise
## 1         2         2
## 2         6         2
## 3         1         6
## 4         7         2
## 5         2        12
## 6         4        45
## 7        11         8
## 8        40         3
## 9         7         2
## 10        0         6
## 11        3         4
## 12        4         0
## 13        5         1
## 14        2         3
## 15       15         0
```

```
View(df)
head(df)
```

```
##      Netflix Exercise
```

```
## 1      2      2
## 2      6      2
## 3      1      6
## 4      7      2
## 5      2     12
## 6      4     45
```

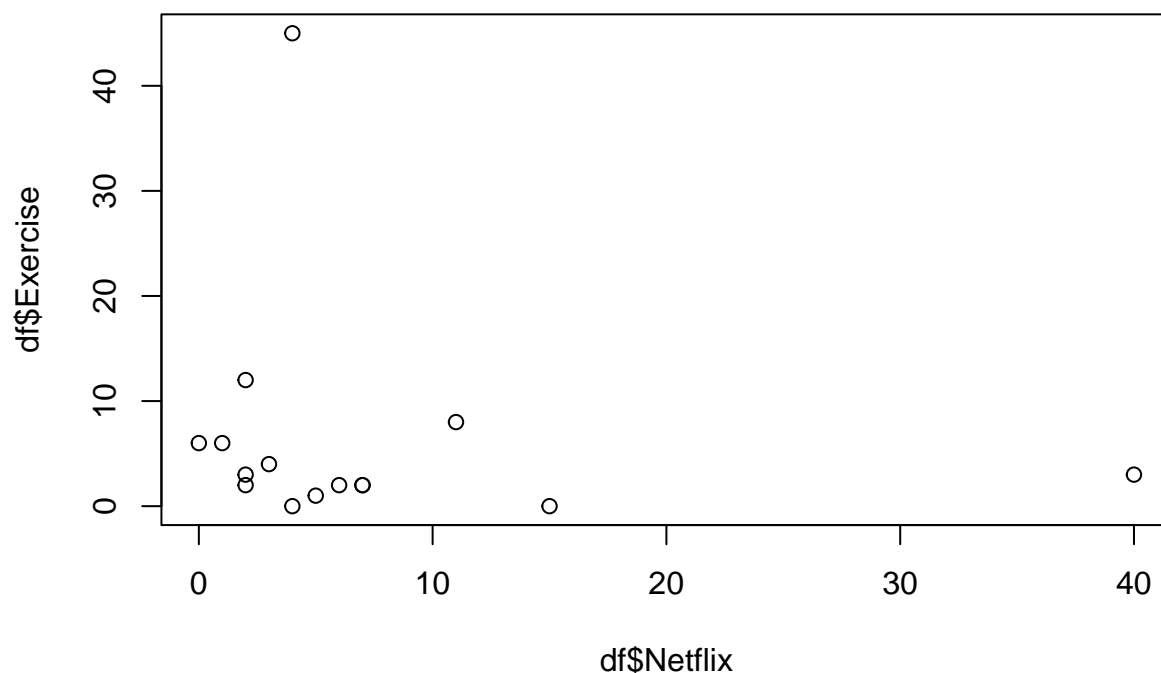
```
tail(df)
```

```
##      Netflix Exercise
## 10         0         6
## 11         3         4
## 12         4         0
## 13         5         1
## 14         2         3
## 15        15         0
```

It's also a good idea to plot your data. This helps you get a general idea for what your data looks like, and to see if there is anything weird going on (i.e., are there any large values in your dataset that could be outliers?)

To make a scatterplot, we can use `plot()` function in R `plot()` works by plotting the first argument on the x-axis, and the second on the y-axis.

```
plot(df$Netflix, df$Exercise)
```



R just automatically sets the axis labels to whatever the input is. However, that's not very pretty, and depending on what you call your columns, it might not be very informative for other people to read. Luckily, we can change the axis labels, and even give our plot a title,

```
plot(df$Netflix, df$Exercise, xlab = "Hours Spent Watching Netflix", #changes the x-axis label
     ylab = "Hours Spent Exercising", # changes the y-axis label
     main = "Plot of Time Spent Watching Netflix vs. Exercising over Break") # gives your plot a title
```

Plot of Time Spent Watching Netflix vs. Exercising over Break



This plot can also give us a general idea of the relation between our variables. For example, here it seems like the more time people spend watching Netflix, the less time they spend exercising. What if we wanted to quantify this relation? We can use the covariance and correlation!

Let's look at covariance first. The covariance between two variables is a measure of how the two variables change together. It only makes sense if there's some connection between the two variables. In this case, each row came from the same person.

The covariance is a bit like the variance, but instead of squared differences from the mean, we multiply these differences from the mean by each other. Let's use the variables in our new dataframe, `df`. Here are the steps:

1. Get the differences from the mean for each variable:

```
diff_nfx <- df$Netflix - mean(df$Netflix, na.rm = TRUE)
diff_ex <- df$Exercise - mean(df$Exercise, na.rm = TRUE)
# you can check this is right -- the differences will sum to 0
sum(diff_ex)
```

```
## [1] -3.552714e-15
```

2. Multiply them by each other (by row – R does this but you can check by hand, or compare the objects to see it)

```
mult_diffs <- diff_nfx * diff_ex
```

3. Sum all these multiplied differences

```
sum_diffs <- sum(mult_diffs, na.rm = TRUE)
```

4. Divide by $N - 1$

```
cov_NetEx <- sum_diffs/14
```

You can verify this yourself by using the `cov()` function

```
cov_NetEx <- cov(df$Netflix, df$Exercise, use = "complete.obs")
cov_NetEx
```

```
## [1] -15.18571
```

For `cov()`, the `use = "complete.obs"` argument acts similarly to `na.rm = T` basically, it will only use data from people who gave an answer to both Netflix and exercise.

We see that the covariance is negative, indicating that the Netflix and Exercise variables are inversely related to each other: higher values for Netflix tend to go with lower values for Exercise, and vice versa!

But how strong is this association? We don't know, because covariances have arbitrary scales based on the scales of the original variables. We don't know how big they could get so we don't know if this value is large or small. That's not very helpful for us, so we need something we know the scale of: correlations.

Correlations can only range between -1 and 1, so they're easier to interpret. We'll standardize the covariance to get a correlation. Standardizing in this case means dividing by the variables' standard deviations. We already know how to get that:

```
sd_N <- sd(df$Netflix, na.rm = T)
sd_E <- sd(df$Exercise, na.rm = T)
```

We divide the covariance by the product of the variances to get a correlation:

```
cov_NetEx/(sd_N*sd_E)
```

```
## [1] -0.1377893
```

And we can check this using the `cor()` function:

```
cor(df$Netflix, df$Exercise, use = "complete.obs")
```

```
## [1] -0.1377893
```

The answers match! Now, since the correlation ranges between -1 and 1, we can say something about the strength of this relation. This value is close to 0, so based on some rules of thumb we can say there is a weak negative relation between watching Netflix and exercise. Of course, what is considered strong vs. weak can depend on the area of research you're in. Next week, we will learn how to conduct a formal statistical test to see whether this correlation is significantly different from 0 or not!