

# Assignment 6

Landon Wilson

2024-11-08

```
library(tidyverse)
library(forcats)
library(nycflights13)
```

## Exercise 1

A common task is to take a set of data that has multiple categorical variables and create a table of the number of cases for each combination. An introductory statistics textbook contains a data set summarizing student surveys from several sections of an intro class. The two variables of interest are **Gender** and **Year** which are the students gender and year in college. *Note: you will need to refer to Chapter 4 and Chapter 7 for some of the operations needed below - this is a great time to review chapter 4!*

a) Download the data set using the following:

```
Survey <- read.csv('https://www.lock5stat.com/datasets2e/StudentSurvey.csv', na.strings=c('',' '))
```

b) Select the specific columns of interest **Year** and **Gender**

```
Survey.2 <- Survey %>%
  select(Year, Gender)

head(Survey.2)
```

```
##      Year Gender
## 1   Senior     M
## 2 Sophomore    F
## 3 FirstYear    M
## 4   Junior     M
## 5 Sophomore    F
## 6 Sophomore    F
```

c) Convert the **Year** column to factors and properly order the factors based on common US progression (FirstYear - Sophomore - Junior - Senior)

```
Survey.3 <- Survey.2 %>%
  drop_na() %>%
  mutate(Year = factor(Year)) %>%
  mutate(Year = fct_relevel(Year, "FirstYear", "Sophomore", "Junior", "Senior"))

levels(Survey.3$Year)
```

```
## [1] "FirstYear" "Sophomore" "Junior" "Senior"
```

d) Convert the **Gender** column to factors and rename them Male/Female.

```
Survey.4 <- Survey.3 %>%
  mutate(Gender = factor(Survey.3$Gender)) %>%
  mutate(Gender = fct_recode(Gender, "Male" = "M"),
         Gender = fct_recode(Gender, "Female" = "F"))

head(Survey.4)
```

```
##      Year Gender
## 1   Senior   Male
## 2 Sophomore Female
## 3 FirstYear   Male
## 4   Junior   Male
## 5 Sophomore Female
## 6 Sophomore Female
```

e) Produce a data set with eight rows and three columns that contains the number of responses for each gender:year combination. *You might want to look at the following functions: `dplyr::count` and `dplyr::drop_na`.*

```
Survey.5 <- Survey.4 %>%
  count(Gender, Year)

Survey.5
```

```
##   Gender      Year  n
## 1 Female FirstYear 43
## 2 Female Sophomore 96
## 3 Female   Junior 18
## 4 Female   Senior 10
## 5  Male FirstYear 51
## 6  Male Sophomore 99
## 7  Male   Junior 17
## 8  Male   Senior 26
```

f) Pivot the table in part (e) to produce a table of the number of responses in the following form:

| Gender | First Year | Sophomore | Junior | Senior |
|--------|------------|-----------|--------|--------|
| Female |            |           |        |        |
| Male   |            |           |        |        |

```
Survey.6 <- Survey.5 %>%
  pivot_wider(names_from = Year,
              values_from = n)

Survey.6
```

```
## # A tibble: 2 x 5
##   Gender FirstYear Sophomore Junior Senior
##   <fct>      <int>      <int>  <int>  <int>
## 1 Female      43        96     18    10
## 2 Male       51        99     17    26
```

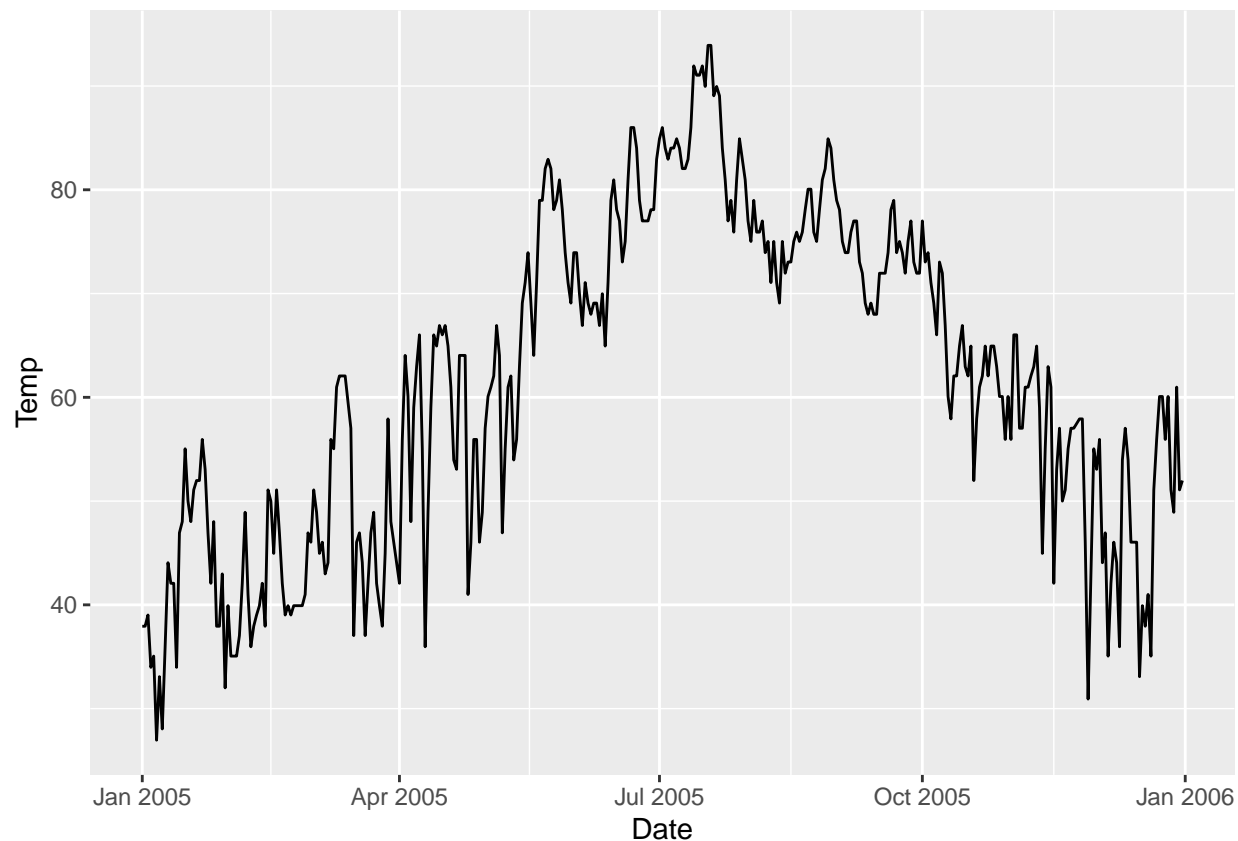
## Exercise 2

From this book's GitHub there is a .csv file of the daily maximum temperature in Flagstaff at the Pulliam Airport. The link is: [https://raw.githubusercontent.com/BuscagliaR/STA\\_444\\_v2/master/data-raw/FlagMaxTemp.csv](https://raw.githubusercontent.com/BuscagliaR/STA_444_v2/master/data-raw/FlagMaxTemp.csv)

a) Create a line graph that gives the daily maximum temperature for 2005. *Make sure the x-axis is a date and covers the whole year.*

```
FlagTemp <- read.csv("https://raw.githubusercontent.com/BuscagliaR/STA_444_v2/master/data-raw/FlagMaxTemp.csv")
FlagTemp.2 <- FlagTemp %>%
  filter(Year == "2005") %>%
  select(!X) %>%
  pivot_longer(cols = starts_with("X"), values_to = "Temp", names_to = "Days") %>%
  mutate(Day = str_extract(Days, "\\d+")) %>%
  mutate(Date = make_date(year = Year, month = Month, day = Day)) %>%
  drop_na() %>%
  select(Date, Temp)

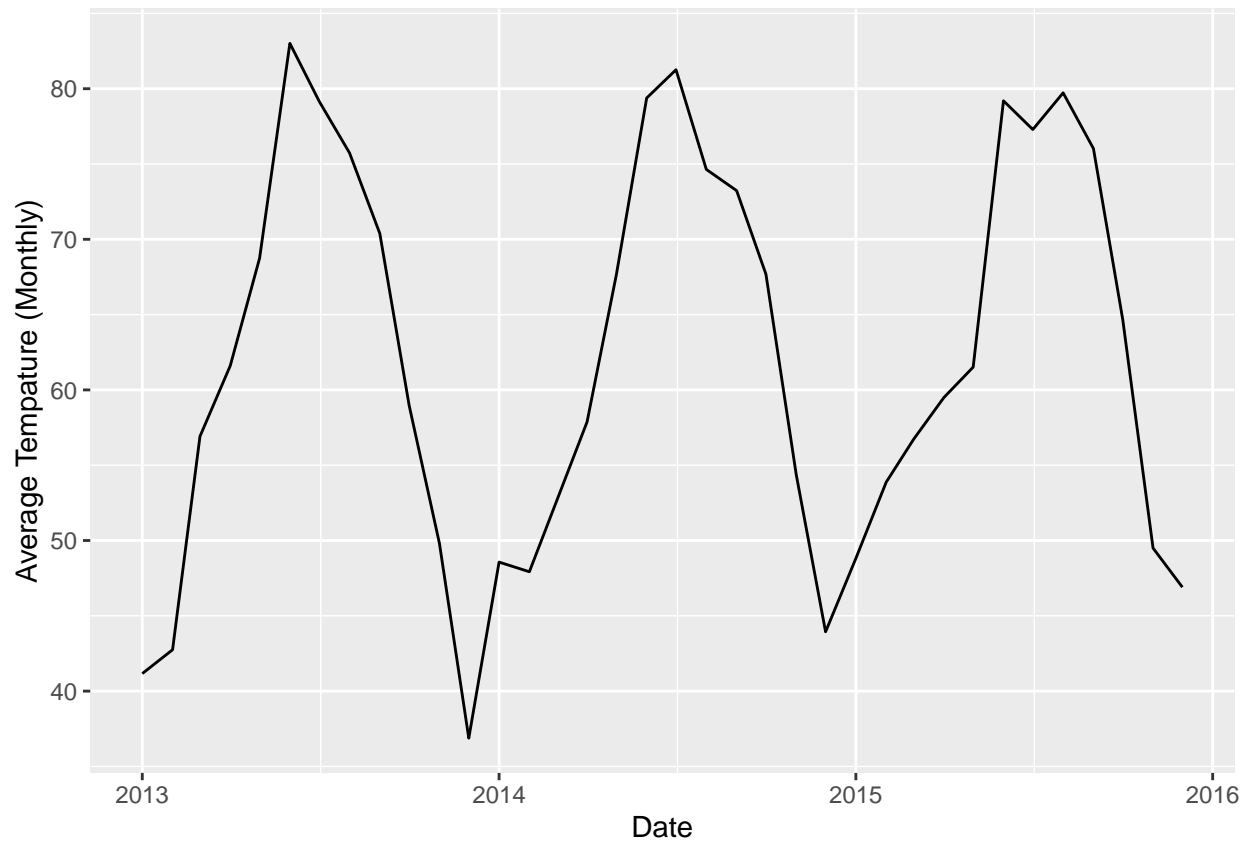
ggplot(FlagTemp.2, aes(x = Date, y = Temp)) +
  geom_line()
```



b) Create a line graph that gives the monthly average maximum temperature for 2013 - 2015. Again the x-axis should be the date and span 3 years.

```
FlagTemp.3 <- FlagTemp %>%
  subset(Year >= "2013" & Year <= "2015") %>%
  select(!X) %>%
  pivot_longer(cols = starts_with("X"), values_to = "Temp", names_to = "Days") %>%
  mutate(Day = str_extract(Days, "\\d+")) %>%
  mutate(Date = make_date(year = Year, month = Month, day = Day)) %>%
  drop_na() %>%
  select(Year, Month, Temp) %>%
  group_by(Year, Month) %>%
  summarise(Avg.Temp = mean(Temp)) %>%
  mutate(Date = make_date(year = Year, month = Month, 1)) %>%
  select(Date, Avg.Temp, Year)

ggplot(FlagTemp.3, aes(x = Date, y = Avg.Temp)) +
  geom_line() +
  labs(y = "Average Temperature (Monthly)")
```



### Exercise 3

For this problem we will consider two simple data sets.

```
A <- tribble(
  ~Name, ~Car,
  'Alice', 'Ford F150',
  'Bob', 'Tesla Model III',
  'Charlie', 'VW Bug')

B <- tribble(
  ~First.Name, ~Pet,
  'Bob', 'Cat',
  'Charlie', 'Dog',
  'Alice', 'Rabbit')
```

a) Combine the data frames together to generate a data set with three rows and three columns using join commands.

```
B <- B %>%
  rename(Name = First.Name)

full_join(A,B)
```

```
## Joining with 'by = join_by(Name)'
```

```
## # A tibble: 3 x 3
##   Name      Car      Pet
##   <chr>   <chr>   <chr>
## 1 Alice   Ford F150   Rabbit
## 2 Bob     Tesla Model III Cat
## 3 Charlie VW Bug      Dog
```

b) It turns out that Alice also has a pet guinea pig. Add another row to the B data set. Do this using either the base function `rbind`, or either of the `dplyr` functions `add_row` or `bind_rows`.

```
new.row <- tibble(Name = "Alice", Pet = "Guinea Pig")

B <- B %>%
  rbind(new.row)
B
```

```
## # A tibble: 4 x 2
##   Name      Pet
##   <chr>   <chr>
## 1 Bob     Cat
## 2 Charlie Dog
## 3 Alice   Rabbit
## 4 Alice   Guinea Pig
```

c) Combine again the A and B data sets together to generate a data set with four rows and three columns using `join` commands.

*Note: You may want to also try using `cbind` to address questions (a) and (c). Leave this as a challenge question and focus on the easier to use `join` functions introduced in this chapter.*

```
full_join(A,B)
```

```
## Joining with 'by = join_by(Name)'
```

```
## # A tibble: 4 x 3
##   Name      Car      Pet
##   <chr>   <chr>   <chr>
## 1 Alice   Ford F150   Rabbit
## 2 Alice   Ford F150   Guinea Pig
## 3 Bob     Tesla Model III Cat
## 4 Charlie VW Bug      Dog
```

## Exercise 4

The package `nycflights13` contains information about all the flights that arrived in or left from New York City in 2013. This package contains five data tables, but there are three data tables we will work with. The data table `flights` gives information about a particular flight, `airports` gives information about a particular airport, and `airlines` gives information about each airline. Create a table of all the flights on February 14th by Virgin America that has columns for the carrier, destination, departure time, and flight duration. Join this table with the airports information for the destination. Notice that because the column for the destination airport code doesn't match up between `flights` and `airports`, you'll have to use the `by=c("TableA.Col"="TableB.Col")` argument where you insert the correct names for `TableA.Col` and `TableB.Col`.

```

data("flights", "airports", "airlines")

flights <- flights %>%
  filter(month == 2, day == 14, carrier == "VX") %>%
  select(carrier, dest, dep_time, air_time)

combo <- left_join(flights, airports, by = c("dest" = "faa"))
combo <- combo %>%
  select(carrier, dest, dep_time, air_time)
combo

```

```

## # A tibble: 10 x 4
##   carrier dest   dep_time air_time
##   <chr>   <chr>   <int>   <dbl>
## 1 VX     LAX       706     347
## 2 VX     SFO       732     344
## 3 VX     LAX       909     341
## 4 VX     LAS       934     307
## 5 VX     SFO     1029     351
## 6 VX     LAX     1317     349
## 7 VX     LAX     1706     335
## 8 VX     SFO     1746     358
## 9 VX     SFO     1852     355
## 10 VX    LAX     2017     337

```