# 01-06: Boxplots 1

## 1 - Purpose

- Order boxplots by factor
- add text to a plot area
- find, and use, quantile values
- put error bars on boxplot

## 2 - Concepts

## 3 - Get data

For this lesson, we will work from the data file we create in the last lesson: LansingNOAA2016-2.csv.

```
1  source( file="scripts/reference.R" );
2  weatherData = read.csv( file="data/LansingNOAA2016-2.csv",
3                          stringsAsFactors = FALSE );
```

## 4 - Creating three new columns

We are going to create, and then use, three new columns in **weatherData**:
- Wind Speed (values: high, medium, low)
- Wind Direction (values: north, south, west, east)
- Change in Temperature (a numeric value giving the difference between high temperatures for consecutive days)

### 4.1 - Wind speed

The column **windSpeed** gives the wind speed in **miles/hour** for each day. We will categorize these wind speed values as **high, medium**, and **low** winds, where:
- **low**: bottom 30% (0-0.3 quantile) of wind speed values
- **medium**: 30%-70% (0.3-0.7 quantile) of wind speed values
- **high**: top 30% (0.7-1.0 quantile) of wind speed values

To do this we need to know what the **30th** (0.3) and **70th** (0.7) percentile (quantile) values for **windSpeed** are.

#### 4.1.1 - Finding percentiles (quantiles)

We can find percentile values using **quantile():**

In **quantiles()**, we pass in the vector we want to check (**windSpeed**) and a vector, called **probs**, that gives the quantile values we want to find.

```
windSpeedQuant = quantile(weatherData[,"windSpeed"], probs=c(.30, .70));
```

The above code says that the **0.30** and **0.70** quantile values will be found in *windSpeed* and these two values will be saved to the variable *windSpeedQuant.* The **0.30** and **0.70** quantile values are **6.2** and **10.4** as shown in *fig 1*.
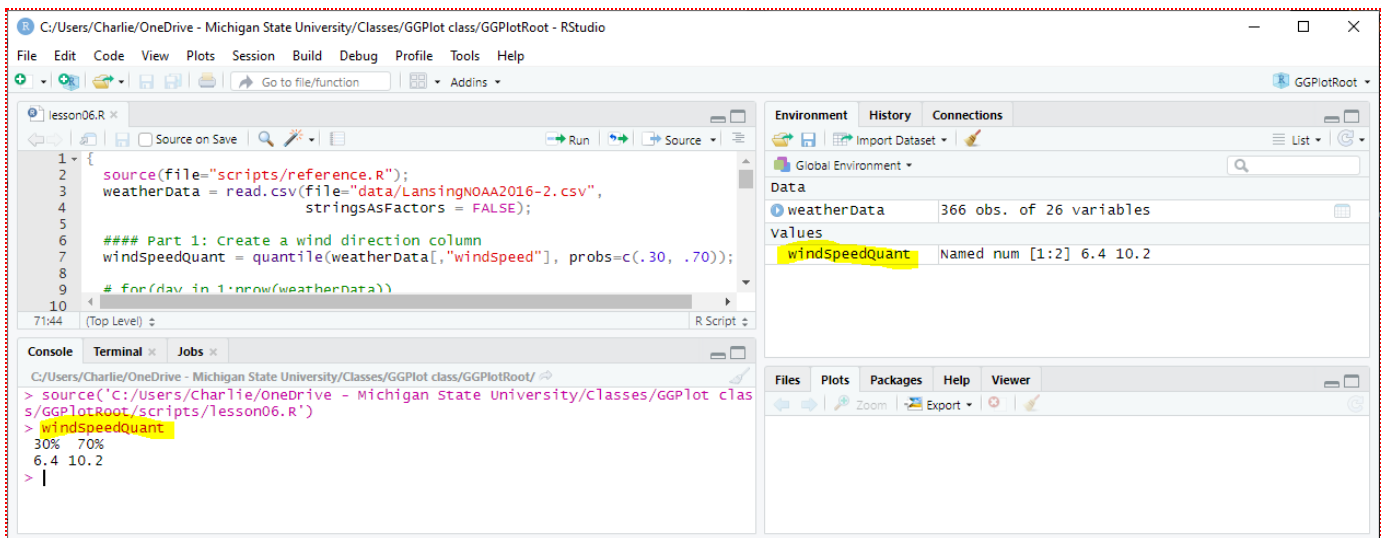


*Fig 1: The 0.30 and 0.70 quantile values for wind speed*

## 4.1.2 - Categorizing using quantiles

So, *windSpeedQuant[1]* = **6.4** and *windSpeedQuant[2]* = **10.2**. Now, we can compare all **366** values in *windSpeed* to the quantile values and place the wind speed values in a category: **low**, **medium**, or **high**.

First, we use a *for()* to go through all rows of *weatherData*. The number of rows is given by *nrow()*

```
for(day in 1:nrow(weatherData))  # nrow(weatherData) = 366
```

*Note: we could have used **for(i in 1:366)** but using **nrow(weatherData)** makes the code more scalable (i.e., it will still work if the number of rows, or days, changes)*

Inside the *for()*, we will use an *if-else* structure to compare each *windSpeed* values to the quantile values in *windSpeedQuant* and, accordingly, set the values in our new column, *windSpeedLevel,* to low, medium, or high.

```
 1  for(day in 1:nrow(weatherData))  # nrow(weatherData) = 366
 2  {
 3     # if the value in windSpeed is less than or equal to the low quant value
 4     if(weatherData[day,"windSpeed"] <= windSpeedQuant[1])       # <= 6.4
 5     {
 6        weatherData[day,"windSpeedLevel"] = "Low";
 7     }
 8     # if the value in windSpeed is greater than or equal to the high quant value
 9     else if(weatherData[day,"windSpeed"] >= windSpeedQuant[2]) # >= 10.2
10     {
11        weatherData[day,"windSpeedLevel"] = "High";
12     }
13     else # the value in windSpeed is between 6.2 and 10.4
14     {
15        weatherData[day,"windSpeedLevel"] = "Medium";
16     }
```
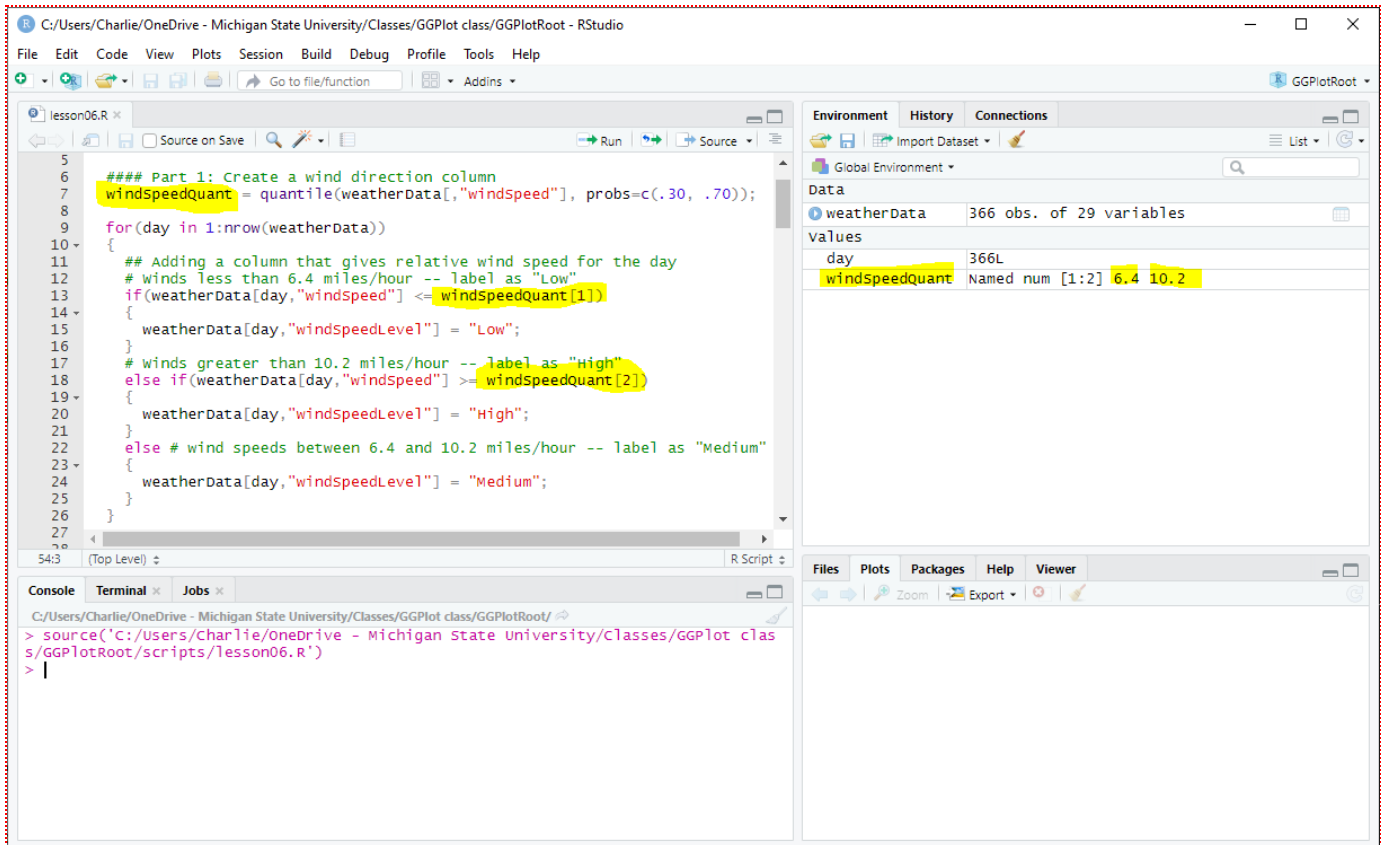
```
    5
    6     #### Part 1: Create a wind direction column
    7     windSpeedQuant = quantile(weatherData[,"windspeed"], probs=c(.30, .70));
    8
    9     for(day in 1:nrow(weatherData))
   10 -  {
   11        ## Adding a column that gives relative wind speed for the day
   12        # Winds less than 6.4 miles/hour -- label as "Low"
   13        if(weatherData[day,"windspeed"] <= windSpeedQuant[1])
   14 -      {
   15          weatherData[day,"windSpeedLevel"] = "Low";
   16        }
   17        # Winds greater than 10.2 miles/hour -- label as "High"
   18        else if(weatherData[day,"windspeed"] >= windSpeedQuant[2])
   19 -      {
   20          weatherData[day,"windSpeedLevel"] = "High";
   21        }
   22        else # wind speeds between 6.4 and 10.2 miles/hour -- label as "Medium"
   23 -      {
   24          weatherData[day,"windSpeedLevel"] = "Medium";
   25        }
   26      }
   27
```

Fig 2: Creating a wind speed level column using quantiles

The new column, **windSpeedLevel**, is shown in *fig 4*.


## 4.2 - Wind direction

The column in **weatherData** called **windSusDir** gives the direction that the most sustained wind came from for the day. The direction is in degrees from **0** to **360** with **0** being due north and **180** being due south.

We are going to convert the degrees to the closest cardinal direction (Fig *3*) and save the value to the column **windDir**:
- **North** is from **315** degrees to **45** degrees (true north is **0** degrees)
- **East** is from **45** degrees to **135** degrees (true East is **90** degrees)
- **South** is from **135** degrees to **225** degrees (true South is **180** degrees)
- **West** is from **225** degrees to **315** degrees (true West is **270** degrees)
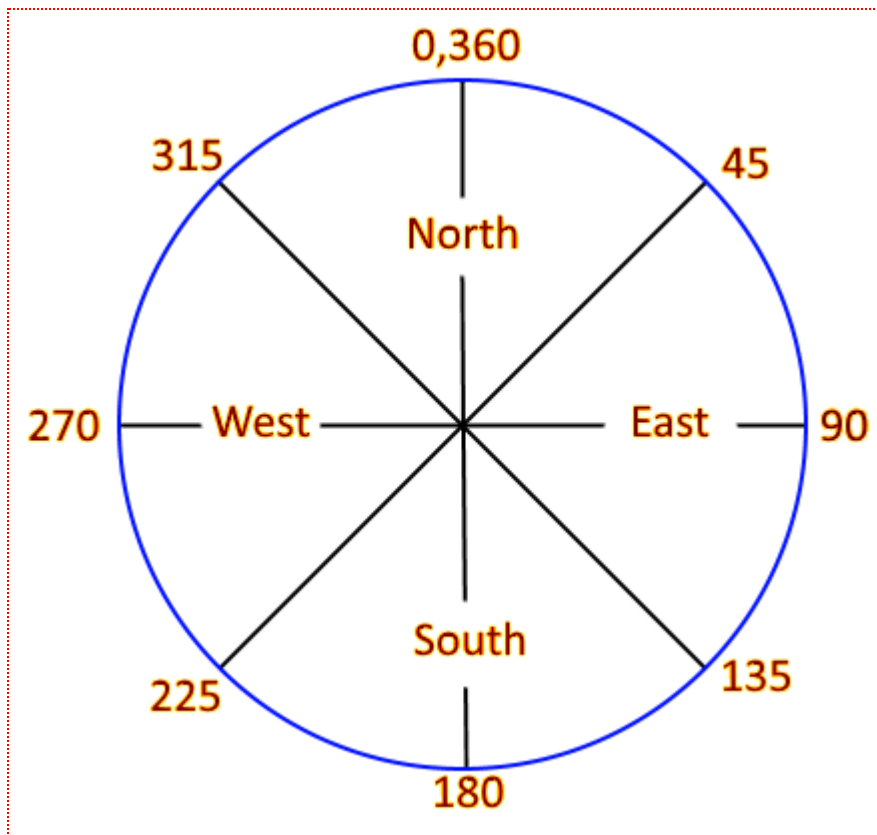
*Fig 3: Compass angles and cardinal directions*

Similar to **windSpeedLevel**, we will use a **for()** to go through every value in the **windSusDir** column and use an **if-else** structure to set a direction for each value in the **windDir** column:

```
1  for(day in 1:nrow(weatherData))
2  {
3    # if the direction is greater than 315 OR less than 45 degrees
4    if(weatherData[day,"windSusDir"] >= 315 ||
5       weatherData[day,"windSusDir"] < 45)
6    {
7      weatherData[day,"windDir"] = "North";
8    }
9    # if the direction is greater than 45 AND less than 135 degrees
10   else if(weatherData[day,"windSusDir"] >= 45 &&
11          weatherData[day,"windSusDir"] < 135)
12   {
13     weatherData[day,"windDir"] = "East";
14   }
15   # if the direction is greater than 135 AND less than 225 degrees
16   else if(weatherData[day,"windSusDir"] >= 135 &&
17          weatherData[day,"windSusDir"] < 225)
18   {
19     weatherData[day,"windDir"] = "South";
20   }
21   else # the direction is between 225 and 315 degrees
22   {
```

```
23        weatherData[day,"windDir"] = "West";
24      }

25 }
```

The new *windDir* column is shown in *fig 4*.

## 4.3 - Change in temperature

Finally, we are going to create a column that gives the difference between the current day's high temperature and the previous day's high temperature.

We need **366** values to complete the whole column, but there is no day previous to the first day on the column.  So, there is no first value and we only have **365** difference values. To fix this issue, we set the first value to *NA*.

Again, we are going to embed an *if-else* structure inside a *for()* that goes through all **366** rows in *weatherData*.  The first value is set to *NA*, the other values are set to the different between the two days.

```
 1 for(day in 1:nrow(weatherData))
 2 {
 3    if( day == 1)  # no day before the first day
 4    {
 5       weatherData[day,"changeMaxTemp"] = NA;
 6    }
 7    else  # subtract previous day's maxTemp from current day's maxTemp
 8    {
 9       weatherData[day,"changeMaxTemp"] = weatherData[day,"maxTemp"] -
10                                          weatherData[day-1,"maxTemp"];
11    }
12 }
```
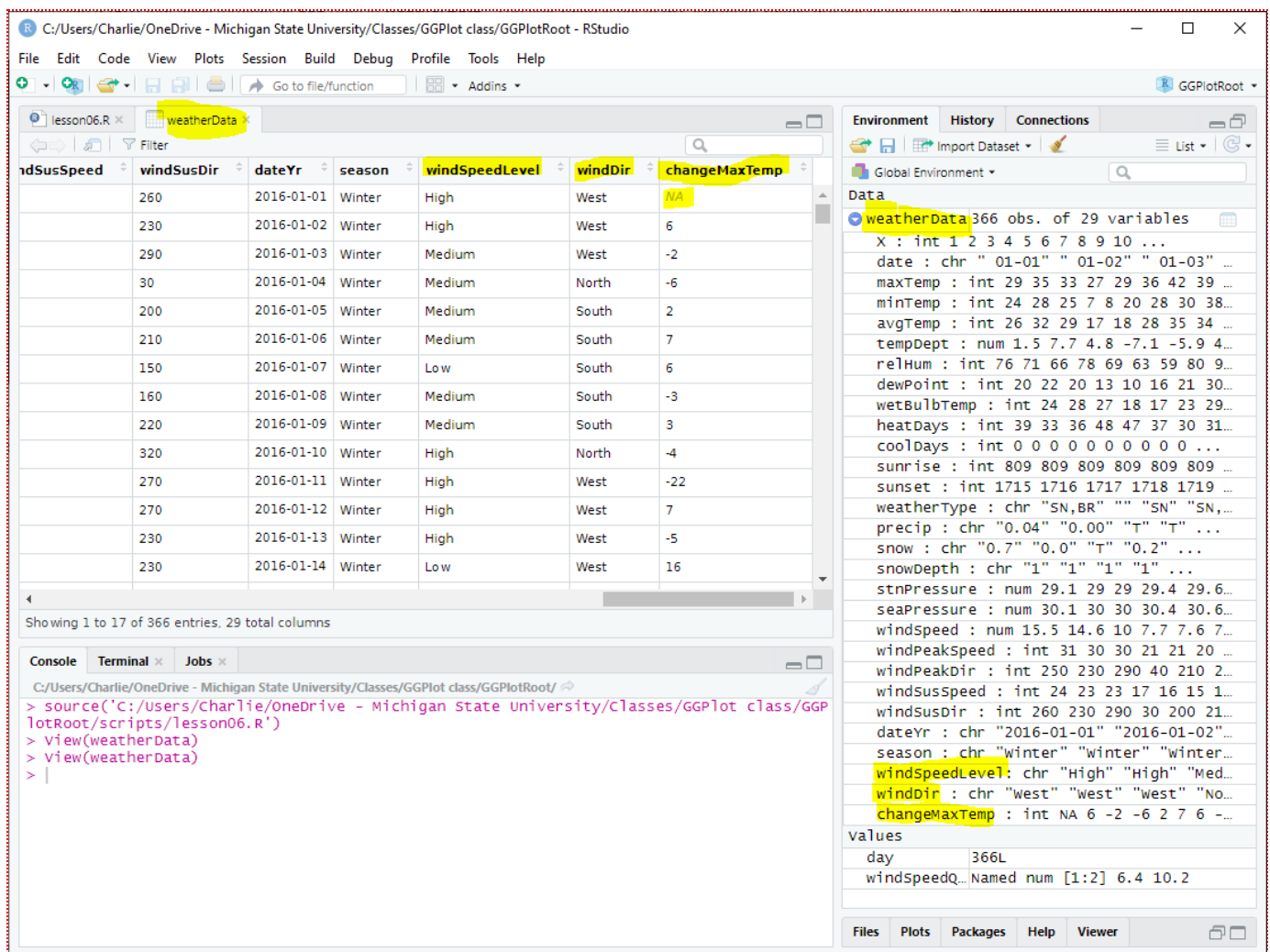
Fig 4: The three new columns in **weatherData**

Extension: one **for()** for all **if-else** structures

## 4.4 - Saving the new dataframe

We are going to use the new dataframe in future lessons, so let's save it to the **data** folder as **LansingNOAA2016-3.csv**:

```
1  write.csv(weatherData, file="data/LansingNOAA2016-3.csv");
```

# 5 - Our first boxplot

The GGPlot component we use to create a boxplot is **geom_boxplot()**. Our first boxplot will use two of the columns we just created: **changeMaxTemp** vs. **windDir**.

```
1  #### Part 2: Plot Wind Speed vs. wind direction
2  thePlot = ggplot(data=weatherData) +
3          geom_boxplot(mapping=aes(x=windDir, y=changeMaxTemp)) +
4          theme_bw() +
5          labs(title = "Change in Temperature vs. Wind Direction",
```

```
6                subtitle = "Lansing, Michigan: 2016",
7                x = "Wind Direction",
8                y = "Degrees (Fahrenheit)");
9  plot(thePlot);
```

In *fig 5* we see that North winds seem to be associated with decreasing temperatures, and South winds seem to be associated with increasing temperatures.  This makes sense as northerly winds often come from the Arctic and southerly winds often come from the Gulf of Mexico.
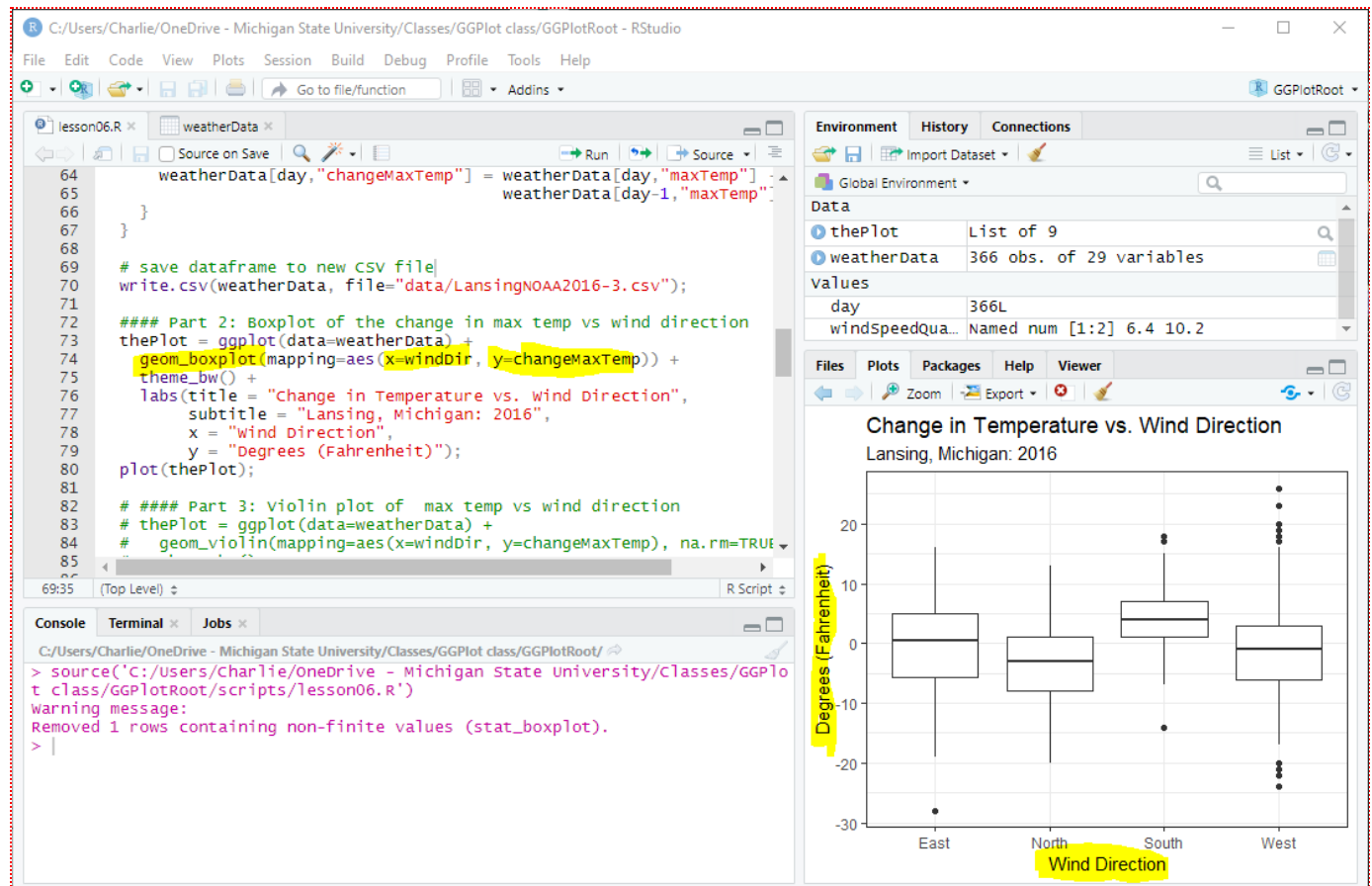


*Fig 5: Our first boxplot: change in temperature vs. wind direction*

*Note: The warning in fig 5 is appearing because there is an **NA** value in the change in temperature data.*

## 5.1 - Variation: violin plot

One variation of a boxplot is the violin plot, which shows distribution instead of quantiles.  The component **geom_violin()** plots a violin plot but, otherwise, works like a **geom_boxplot().**

The parameter **na.rm=TRUE** was added to stop GGPlot from warning us every time that there is an **NA** value in the data (because the first value in **changeMaxTemp** is **NA**).  This will be added to all subsequent plots.

```
1  #### Part 3: Violin plot of Wind Speed vs. wind direction
2  thePlot = ggplot(data=weatherData) +
3          geom_violin(mapping=aes(x=windDir, y=changeMaxTemp),
4                  na.rm=TRUE) +
```

```
 5            theme_bw() +
 6            labs(title = "Change in Temperature vs. Wind Direction",
 7                 subtitle = "Lansing, Michigan: 2016",
 8                 x = "Wind Direction",
 9                 y = "Degrees (Fahrenheit)");
10 plot(thePlot);
```
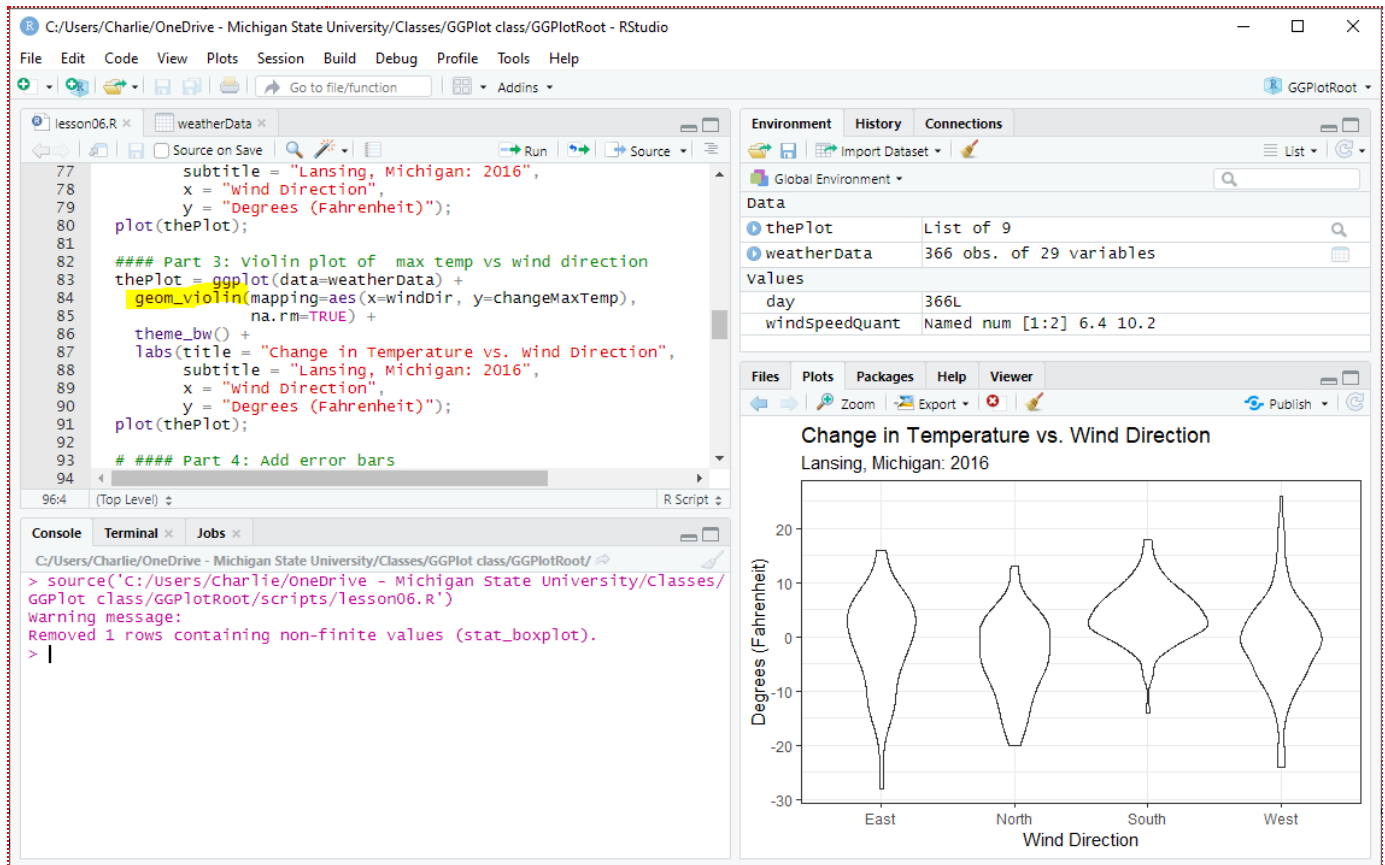


*Fig 6: Violin plot of change in temperature vs wind direction*

## 5.2 - Variation: Bars at the end of the whiskers

**geom_boxplot()** has no support for bars on the end of the boxplot whiskers. One somewhat hacky way to add bars is to create another plotting component on top of the boxplot. The component is **stat_boxplot()**, which contains error bars.

In the following code, **stat_boxplot()** plots just the whiskers and error bars. The parameter **width** takes multipliers *values* with **1** meaning the error bar takes up the whole column (i.e., same width as the boxplot). So, **width=0.2** means the error bar takes up 20% of the column.

```
1 #### Part 4: Add error bars
2 thePlot = ggplot(data=weatherData) +
3         stat_boxplot(mapping=aes(x=windDir, y=changeMaxTemp),
4                 na.rm=TRUE,
5                 geom = "errorbar",
6                 width = 0.2) +
7         geom.boxplot(mapping.aes(x=windDir, y=changeMaxTemp)
```

```
 7            geom_boxplot(mapping=aes(x=windDir, y=changeMaxTemp),
 8                    na.rm=TRUE) +
 9         theme_bw() +
10         labs(title = "Change in Temperature vs. Wind Direction",
11             subtitle = "Lansing, Michigan: 2016",
12             x = "Wind Direction",
13             y = "Degrees (Fahrenheit)");
14 plot(thePlot);
```

*Note: **stat_boxplot()** is coded before and, hence, plotted behind **geom_boxplot()***
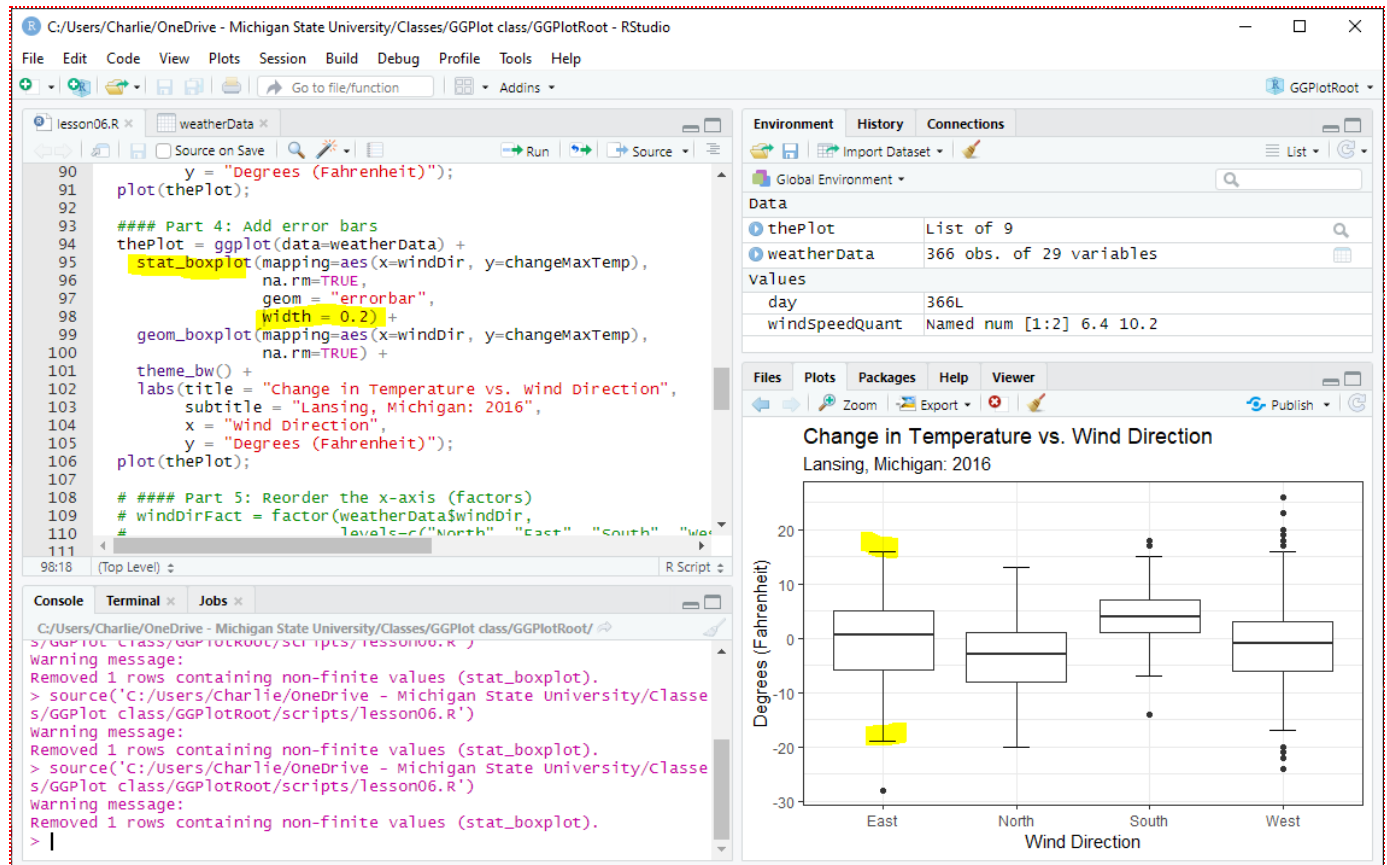


Fig 7: Adding error bars to a boxplot (there is currently no elegant solution for this)

# 6 - Reordering x-axis values

The x-axis (Wind Direction) labels go in alphabetical order: **East**, **North**, **South**, **West**.  GGPlot will default to an alphabetical ordering if factor levels are not provided.  We can force the issue by creating a factor, called ***windDirFact***, with the levels provided in the preferred order:

```
1 ### Part 5: Re-order the directions on the x-axis using factor(s)
2 windDirFact = factor(weatherData$windDir,
3             levels=c("North", "East", "South", "West"));
```

Now we use the factor *windDirFact* in place of *windDir* -- just remember to replace *windDir* with *winDirFact* in both the *stat_boxplot()* and *geom_boxplot()* components.

```
 1  thePlot = ggplot(data=weatherData) +
 2          stat_boxplot(mapping=aes(x=windDirFact, y=changeMaxTemp),
 3                      na.rm=TRUE,
 4                      geom = "errorbar",
 5                      width = 0.2) +
 6          geom_boxplot(mapping=aes(x=windDirFact, y=changeMaxTemp),
 7                      na.rm=TRUE) +
 8          theme_bw() +
 9          labs(title = "Change in Temperature vs. Wind Direction",
10              subtitle = "Lansing, Michigan: 2016",
11              x = "Wind Direction",
12              y = "Degrees (Fahrenheit)");
13  plot(thePlot);
```

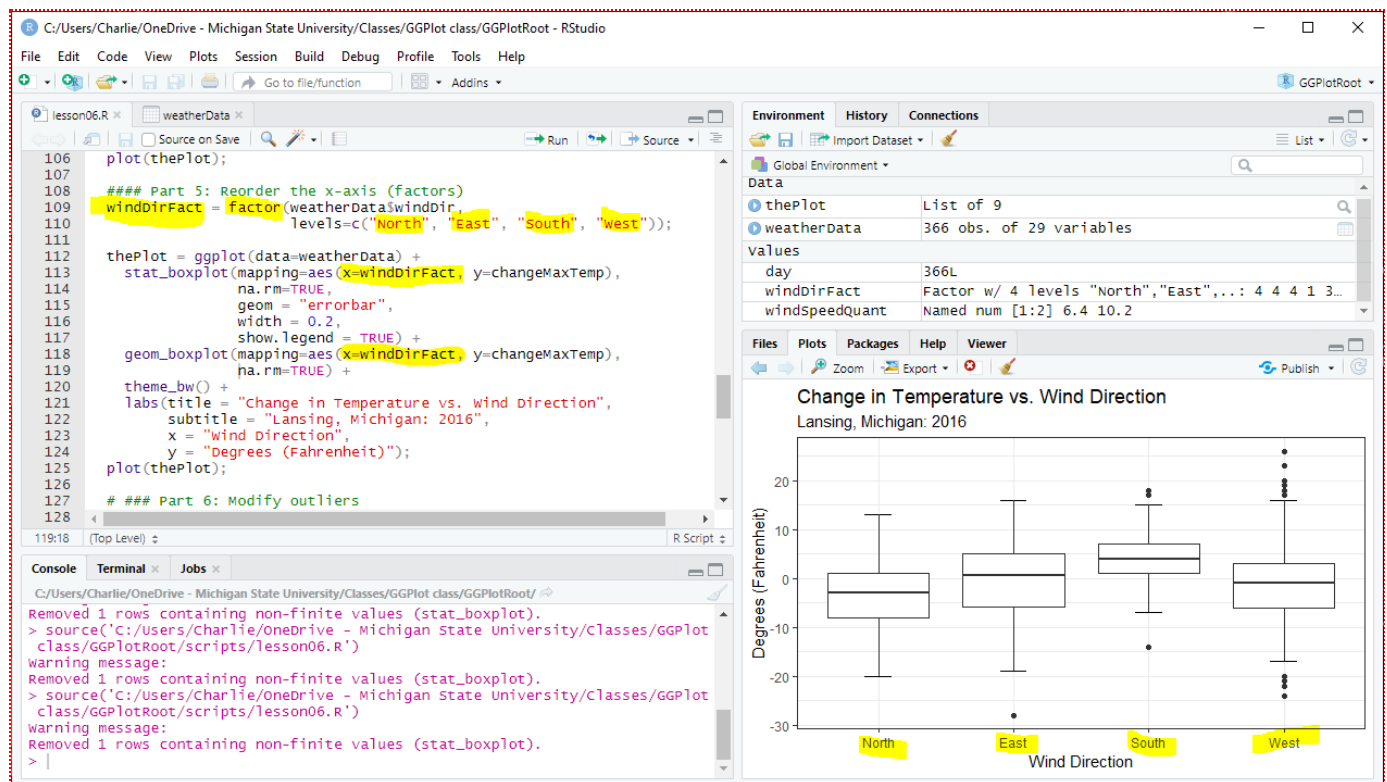The boxplots are rearranged in order of the factor levels.



*Fig 8: Reordering the categorical values on the x-axis*

# 7 - Modifying Outliers

There are many parameters in the component *geom_boxplot()* that relate to outliers:

- *outlier.shape*: the shape number (1-25) or character (in quotes) used to represent the outlier
  *note: outlier.shape = NA will remove the outliers from the plot*
- *outlier.color*: outlier color

- **outlier.size**: outlier size, as a numeric multiplier
- **outlier.alpha**: outlier transparency (from **0**-invisible to **1**-opaque) -- good for situations where outliers overlap
- **outlier.fill**: outlier fill color (only **outlier.shape** between **21-25** use fill color)
- **outlier.stroke**: size of the border as a multiplier (only useful for shapes that have a fill -- otherwise, use **outlier.size**)

In the following example, I change four of the outlier parameters

```
1  thePlot = ggplot(data=weatherData) +
2          stat_boxplot(mapping=aes(x=windDirFact, y=changeMaxTemp),
3                       na.rm=TRUE,
4                       geom = "errorbar",
5                       width = 0.2) +
6          geom_boxplot(mapping=aes(x=windDirFact, y=changeMaxTemp),
7                       na.rm=TRUE,
8                       outlier.shape = "@",
9                       outlier.color = "red",
10                      outlier.alpha = 0.6,
11                      outlier.size = 4 ) +
12         theme_bw() +
13         labs(title = "Change in Temperature vs. Wind Direction",
14             subtitle = "Lansing, Michigan: 2016",
15             x = "Wind Direction",
16             y = "Degrees (Fahrenheit)");
17 plot(thePlot);
```

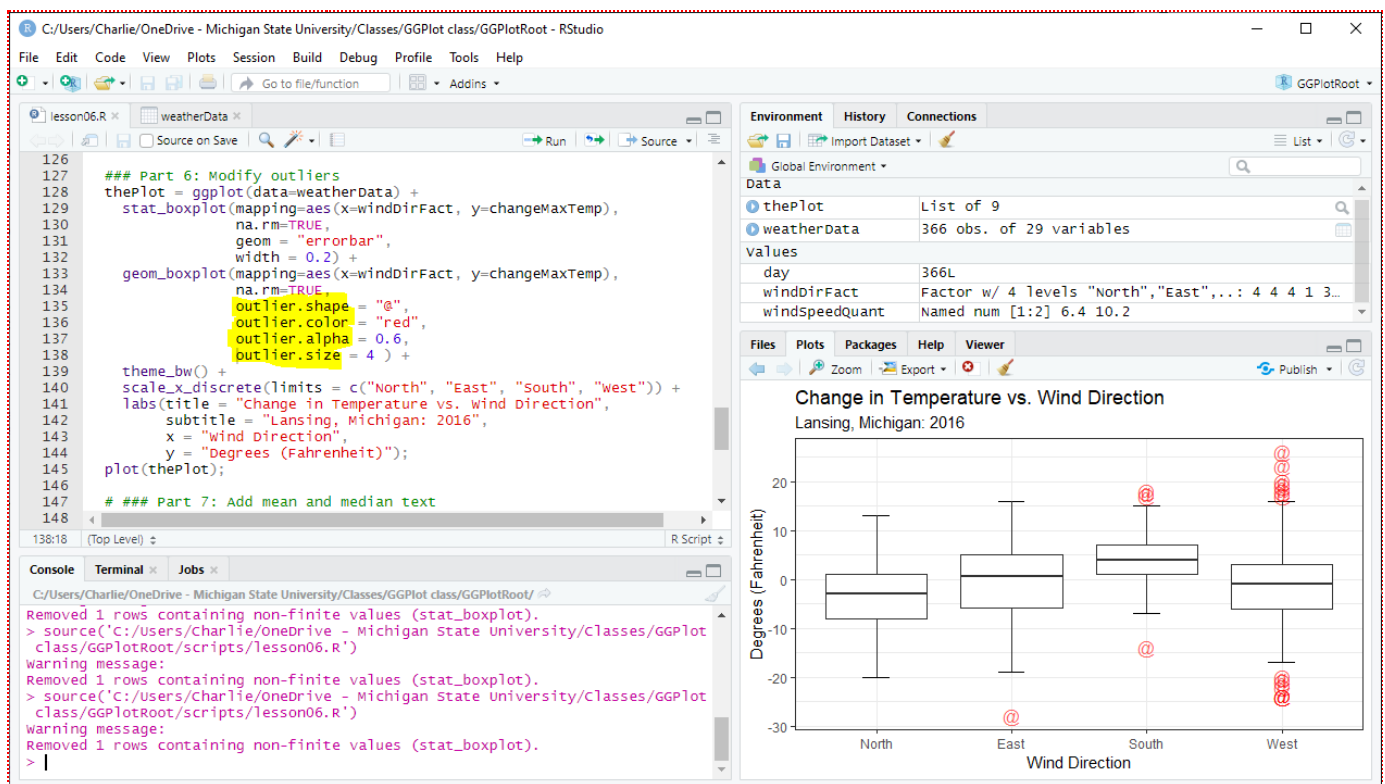The outliers are now red, 60% opaque (40% transparent), 4mm tall, and shaped like an "@".

*Fig 9: Modifying the shape, color, alpha, and size of the outliers*

Note: the author recognizes the tacky nature of the outliers' styling -- this is being used for demo purposes. For the next plot, I will set **outlier.shape** to **NA** to remove the outliers.

# 8 - Adding text to the plot

We are going to add text to (i.e., annotate) the plot.  The text will display the median **Change In Temperature** values for days where the winds come from the **North** and days where the winds come from the **South**.

To this we need to:
- Find the rows in **weatherData** that have **North** winds and the rows in **weatherData** that have **South** winds
- Find the median **changeMaxTemp** for rows that have **North** winds and rows that have **South** winds
- Add the median value to the plot using the **annotate()** component

## 8.1 - Find rows which() meet a condition

We use **which()** to find the indices of a vector that meet a condition.  In this case the condition is: which values in **windDir** are equal to **North**.  The index of the values are stored in **northVals**.

```
northVals=which(weatherData$windDir == "North");
```

And which values in **windDir** are equal to **South**

```
southVals=which(weatherData$windDir == "South");
```

**northVals** and **southVals** are a vector of numbers.  There are **61 North** values (starting with days 4,10, and 23) representing the 61 days that had northerly winds. There are **101 South** values (starting with days 5, 6, and 7) representing the 101 days that had southerly winds (*fig 10*).
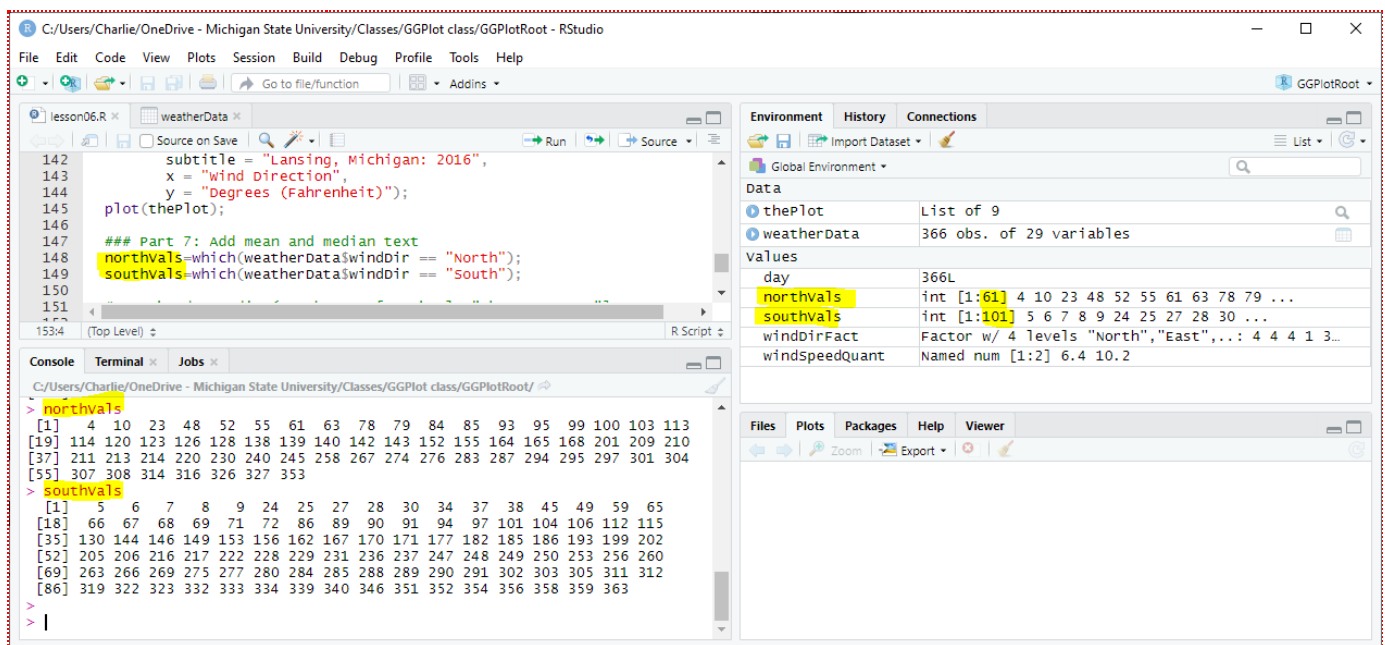
*Fig 10: Finding the values in a vector that meet a condition using **which()***

## 8.2 - Find median value of indexed values

***northVals*** is a vector with all the days (i.e., rows) that had northerly winds. ***southVals*** is a vector with all the days (i.e., rows) that had southerly winds.

We want a vector that contains the change in temperature, ***changeMaxTemp**,* for the 61 days in ***northVals:***

```
weatherData[northVals, "changeMaxTemp"]
```

And a vector that contains ***changeMaxTemp*** values for the 101 days in ***southVals***:

```
weatherData[southVals, "changeMaxTemp"]
```

To get the median of the indexed values above, we use ***median()*** and set ***na.rm = TRUE*** because we want to ignore the **NA** values in ***changeMaxTemp.***

In the following code, the median values for noth and south winds are saved to the variables ***northMed*** and ***southMed***.

```
1  northMed = median(weatherData[northVals,"changeMaxTemp"], na.rm=TRUE);
2  southMed = median(weatherData[southVals,"changeMaxTemp"], na.rm=TRUE);
```
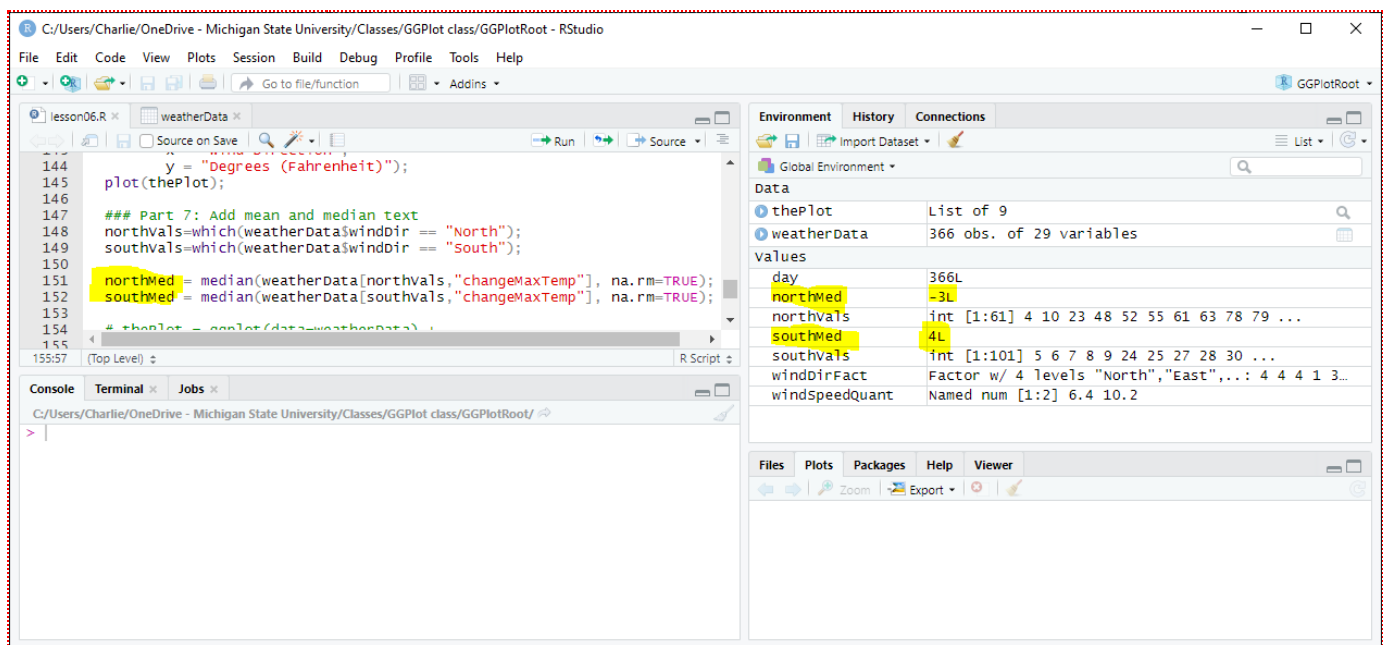
*Fig 11: Median **changeMaxTemp** values for days with **North** and **South** winds*

## 8.3 - Annotate plot with median values

Now we just need to put the median values (***northMed*** and ***southMed***) on the plot using the ***annotate()*** component.

The subcomponents we set in ***annotate()*** are:
- **geom**: the type of value we are putting on the plot (in this case, **text**)
- **x**: x-position of the text -- for discrete axis, x=1 is the location of the first plot, x=2 is the second plot...
- **y**: y-position of the text -- for continuous plot, this value matches the value on the axis
- **color**: color of text
- **label**: the text on the plot -- using ***paste()***, we will put together the text "median:" and the variable that holds the calculated median value

We are adding two ***annotate()*** components added to the plot -- one for **North**, the other for **South**

```
1  thePlot = ggplot(data=weatherData) +
2          stat_boxplot(mapping=aes(x=windDirFact, y=changeMaxTemp),
3                  na.rm=TRUE,
4                  geom = "errorbar",
5                  width = 0.2) +
6          geom_boxplot(mapping=aes(x=windDirFact, y=changeMaxTemp),
7                  na.rm=TRUE,
8                  outlier.shape = "@",
9                  outlier.color = "red",
10                 outlier.alpha = 0.6,
11                 outlier.size = 4 ) +
12         annotate(geom="text",      # North median
13                 x=1,
14                 y=20,
15                 color="blue",
16                 label paste("median:"  northMed) )
```

```
16                  label=paste("median:", northMed) ) +
17            annotate(geom="text",     # South median
18                  x=3,
19                  y=-10,
20                  color="red",
21                  label=paste("median:", southMed) ) +
22            theme_bw() +
23            labs(title = "Change in Temperature vs. Wind Direction",
24                  subtitle = "Lansing, Michigan: 2016",
25                  x = "Wind Direction",
26                  y = "Degrees (Fahrenheit)");
27 plot(thePlot);
```

The text for the **North** median is on column **1** at the **20** degree line.
The text for the **South** median is on column **3** at the **-10** degree line.
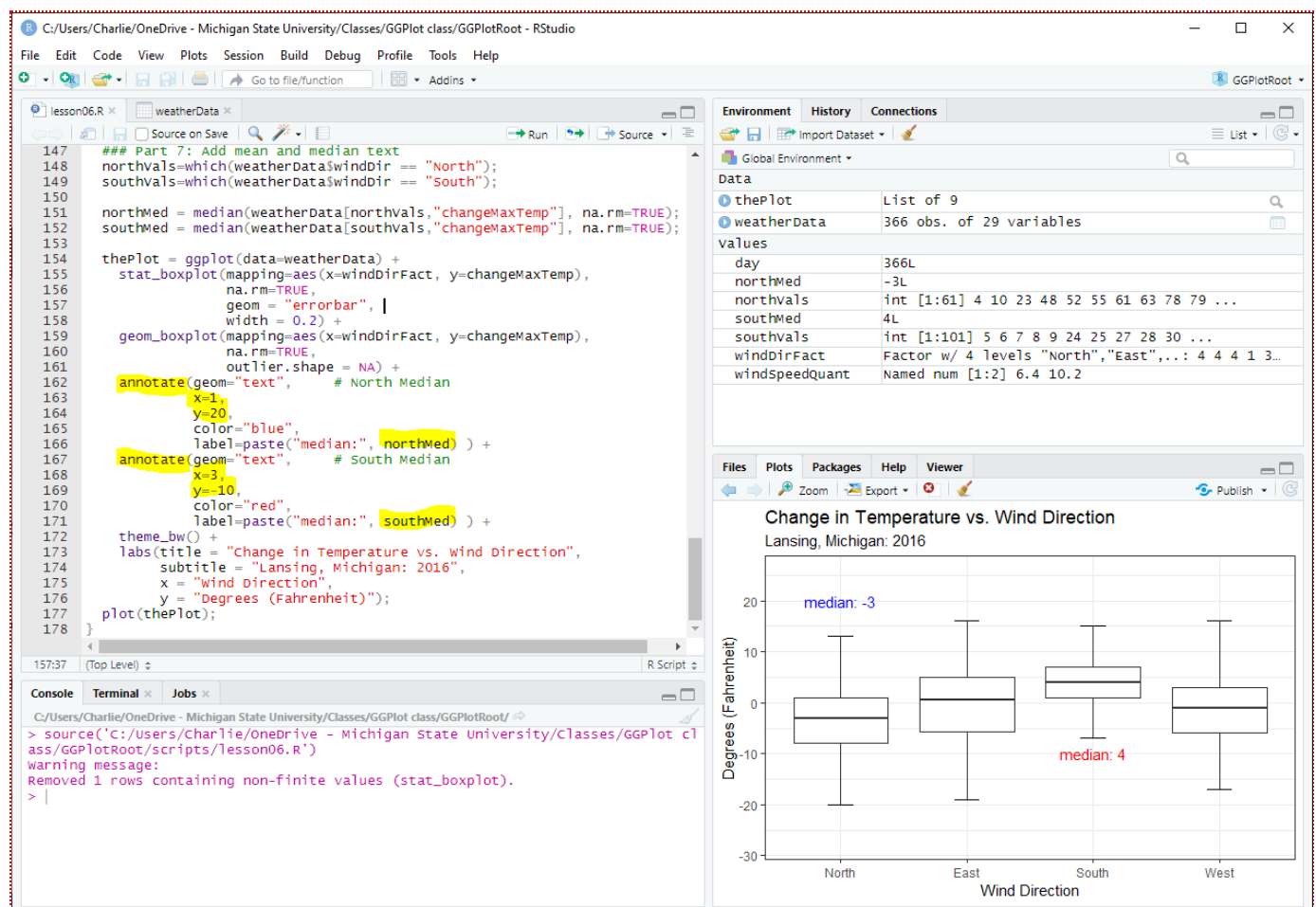


*Fig 12: Adding median values to the plot using **annotate()***

# 9 - Application

1. Create a script file in your GGPlot Class Project called ***app06.r***.
2. Create a vector that holds the quantile values for the 20, 40, 60, and 80[th] percentile of ***stnPressure***.

3. Create a new column called **pressureLevel** that creates five evenly spaced levels for **stnPressure**
4. Make a boxplot of **windSusSpeeds** vs **pressureLevel**
5. Find the dates for the three outliers at the lowest pressure level.  The three outliers are also the three highest **windSusSpeed** for the year.
6. Label the outliers for the lowest pressure level with the dates on the plot.

## 10 - Extension: One for() for all if-else structures

The three **for()** used to create the three new column in **weatherData** (**windSpeedLevel**, **windDir**, and **changeMaxTemp**) could all be condensed into one **for()** because the three **for()** all execute over the same number of iterations, **366**.

```r
 1 for(day in 1:nrow(weatherData))
 2 {
 3   ## Adding a column that gives relative wind speed for the day
 4   # Winds less than 6.4 miles/hour -- label as "Low"
 5   if(weatherData[day,"windSpeed"] <= windSpeedQuant[1])
 6   {
 7     weatherData[day,"windSpeedLevel"] = "Low";
 8   }
 9   # Winds greater than 10.2 miles/hour -- label as "High"
10   else if(weatherData[day,"windSpeed"] >= windSpeedQuant[2])
11   {
12     weatherData[day,"windSpeedLevel"] = "High";
13   }
14   else # wind speeds between 6.4 and 10.2 miles/hour -- label as "Medium"
15   {
16     weatherData[day,"windSpeedLevel"] = "Medium";
17   }
18
19   ## Adding a column that gives the cardinal wind direction
20   # if the direction is greater than 315 OR less than 45 degrees
21   if(weatherData[day,"windSusDir"] >= 315 ||
22   weatherData[day,"windSusDir"] < 45)
23   {
24     weatherData[day,"windDir"] = "North";
25   }
26   # if the direction is greater than 45 AND less than 135 degrees
27   else if(weatherData[day,"windSusDir"] >= 45 &&

28   weatherData[day,"windSusDir"] < 135)
29   {
30     weatherData[day,"windDir"] = "East";
31   }
32   # if the direction is greater than 135 AND less than 225 degrees
33   else if(weatherData[day,"windSusDir"] >= 135 &&
```

```r
34                weatherData[day,"windSusDir"] < 225)
35   {
36     weatherData[day,"windDir"] = "South";
37   }
38   else # the directions is between 225 and 315 degrees
39   {
40     weatherData[day,"windDir"] = "West";
41   }
42   ### Adding a changeMaxTemp column
43   if(day == 1)
44   {
45     weatherData[day,"changeMaxTemp"] = NA;
46   }
47   else
48   {
49     weatherData[day,"changeMaxTemp"] = weatherData[day,"maxTemp"] -
50     weatherData[day-1,"maxTemp"];
51   }
52 }
```