

Intro to Data Science HW 7

```
# Enter your name here: Nora Lin
```

Copyright Jeffrey Stanton, Jeffrey Saltz, and Jasmina Tacheva

Attribution statement: (choose only one and delete the rest)

```
# 3. I did this homework with help from Patrick Furlong but did not cut and paste any code.
```

The chapter on **linear models** (“Lining Up Our Models”) introduces **linear predictive modeling** using the tool known as **multiple regression**. The term “multiple regression” has an odd history, dating back to an early scientific observation of a phenomenon called “**regression to the mean**.” These days, multiple regression is just an interesting name for using **linear modeling** to assess the **connection between one or more predictor variables and an outcome variable**.

In this exercise, you will **predict Ozone air levels from three predictors**.

- A. We will be using the **airquality** data set available in R. Copy it into a dataframe called **air** and use the appropriate functions to **summarize the data**.

```
air <- airquality
summary(air)
```

```
##      Ozone      Solar.R      Wind      Temp
##  Min.   : 1.00   Min.   : 7.0   Min.   : 1.700   Min.   :56.00
## 1st Qu.: 18.00   1st Qu.:115.8   1st Qu.: 7.400   1st Qu.:72.00
## Median : 31.50   Median :205.0   Median : 9.700   Median :79.00
## Mean   : 42.13   Mean   :185.9   Mean   : 9.958   Mean   :77.88
## 3rd Qu.: 63.25   3rd Qu.:258.8   3rd Qu.:11.500   3rd Qu.:85.00
## Max.   :168.00   Max.   :334.0   Max.   :20.700   Max.   :97.00
## NA's   :37      NA's   :7
##      Month      Day
##  Min.   :5.000   Min.   : 1.0
## 1st Qu.:6.000   1st Qu.: 8.0
## Median :7.000   Median :16.0
## Mean   :6.993   Mean   :15.8
## 3rd Qu.:8.000   3rd Qu.:23.0
## Max.   :9.000   Max.   :31.0
##
```

- B. In the analysis that follows, **Ozone** will be considered as the **outcome variable**, and **Solar.R**, **Wind**, and **Temp** as the **predictors**. Add a comment to briefly explain the outcome and predictor variables in the dataframe using **?airquality**.

```
#Predictors variables are independent  
#Outcome variables are dependent
```

```
#From the help function:
```

```
#Ozone: Mean ozone in parts per billion from 1300 to 1500 hours at Roosevelt Island
```

```
#Solar.R: Solar radiation in Langleys in the frequency band 4000-7700 Angstroms from 0800 to 1200 hours
```

```
#Wind: Average wind speed in miles per hour at 0700 and 1000 hours at LaGuardia Airport
```

```
#Temp: Maximum daily temperature in degrees Fahrenheit at La Guardia Airport.
```

C. Inspect the outcome and predictor variables – are there any missing values? Show the code you used to check for that.

```
air$Ozone[is.na(air$Ozone)]
```

```
## [1] NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA  
## [26] NA NA NA NA NA NA NA NA NA NA NA NA NA NA
```

```
air$Solar.R[is.na(air$Solar.R)]
```

```
## [1] NA NA NA NA NA NA NA NA
```

```
air$Wind[is.na(air$Wind)]
```

```
## numeric(0)
```

```
air$Temp[is.na(air$Temp)]
```

```
## integer(0)
```

```
#Ozone has 37 missing values
```

```
#Solar.R has 7 missing values
```

```
#Wind and temp have no missing values
```

D. Use the `na_interpolation()` function from the **imputeTS** package (remember this was used in a previous HW) to fill in the missing values in each of the 4 columns. Make sure there are no more missing values using the commands from Step C.

```
#install.packages("imputeTS")  
library(imputeTS)
```

```
## Registered S3 method overwritten by 'quantmod':  
## method from  
## as.zoo.data.frame zoo
```

```
air$Ozone <- na_interpolation(air$Ozone)
air$Solar.R <- na_interpolation(air$Solar.R)
```

- E. Create **3 bivariate scatterplots (X-Y) plots** (using ggplot), for each of the predictors with the outcome. **Hint:** In each case, put **Ozone on the Y-axis**, and a **predictor on the X-axis**. Add a comment to each, describing the plot and explaining whether there appears to be a **linear relationship** between the outcome variable and the respective predictor.

```
#install.packages("tidyverse")
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.5      v purrr  0.3.4
## v tibble  3.1.5      v dplyr  1.0.7
## v tidyr   1.1.4      v stringr 1.4.0
## v readr   2.0.2      v forcats 0.5.1
```

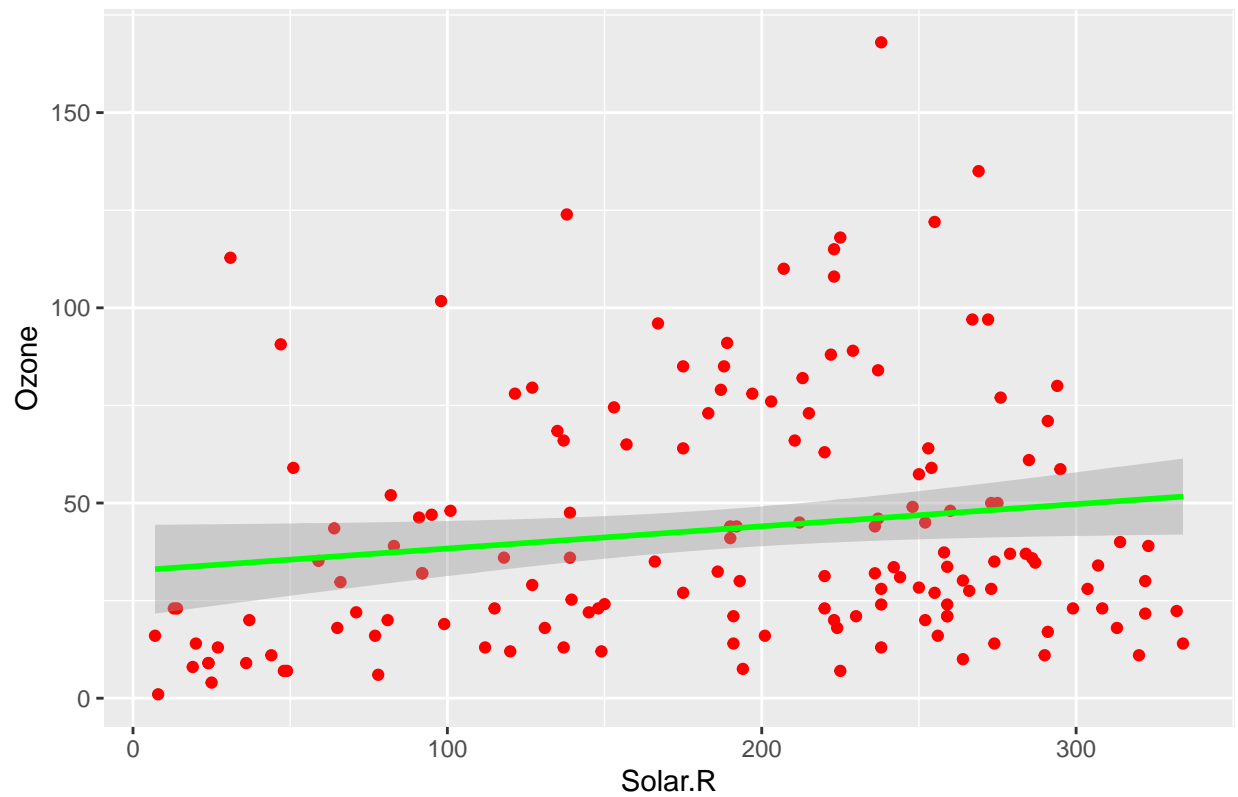
```
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
#install.packages("ggplot2")
library(ggplot2)
```

```
#Scatterplot using Ozone as outcome and Solar.R as predictor
ggplot(air)+
  geom_point(aes(Solar.R,Ozone),color='red')+
  stat_smooth(aes(Solar.R,Ozone),method='lm',color='green')+
  ggtitle("Scatterplot of Ozone by Solar.R")
```

```
## 'geom_smooth()' using formula 'y ~ x'
```

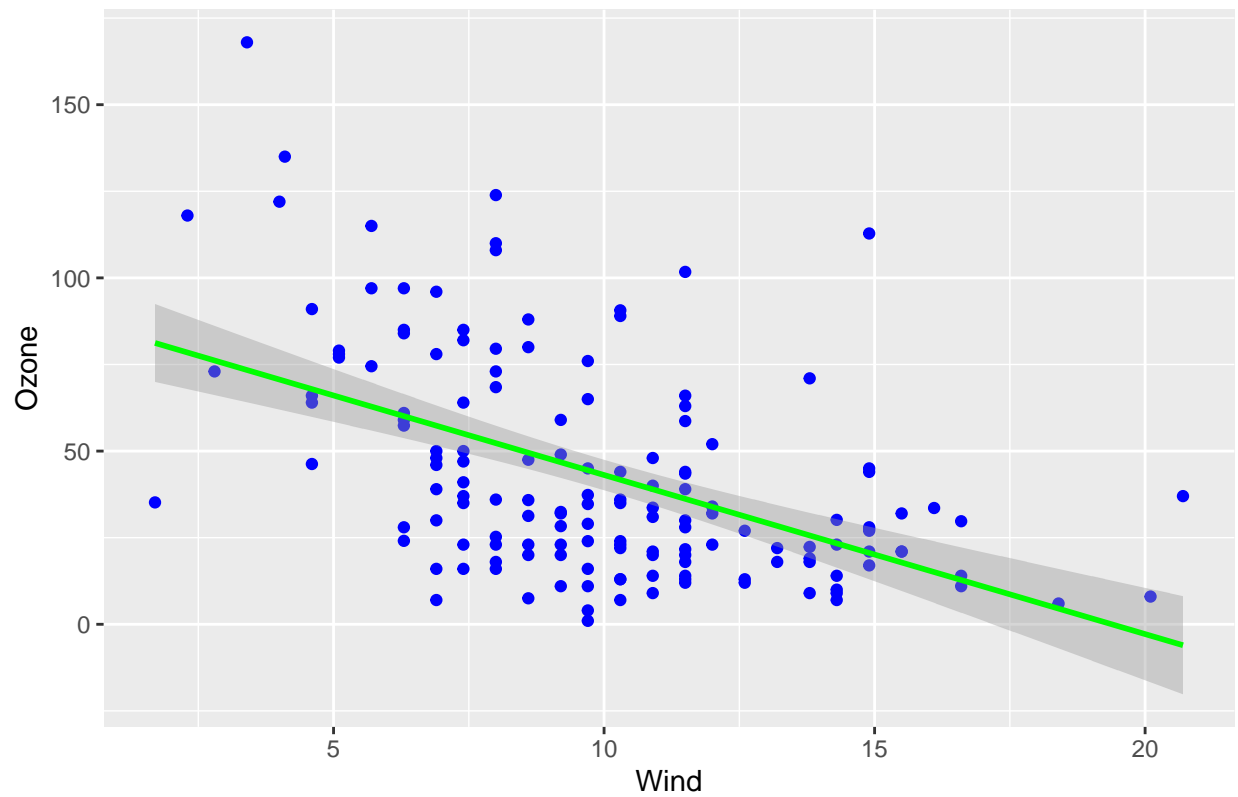
Scatterplot of Ozone by Solar.R



```
#Scatterplot using Ozone as outcome and Wind as predictor  
ggplot(air)+  
  geom_point(aes(Wind,Ozone),color='blue')+  
  stat_smooth(aes(Wind,Ozone),method='lm',color='green')+  
  ggtitle("Scatterplot of Ozone by Wind")
```

```
## 'geom_smooth()' using formula 'y ~ x'
```

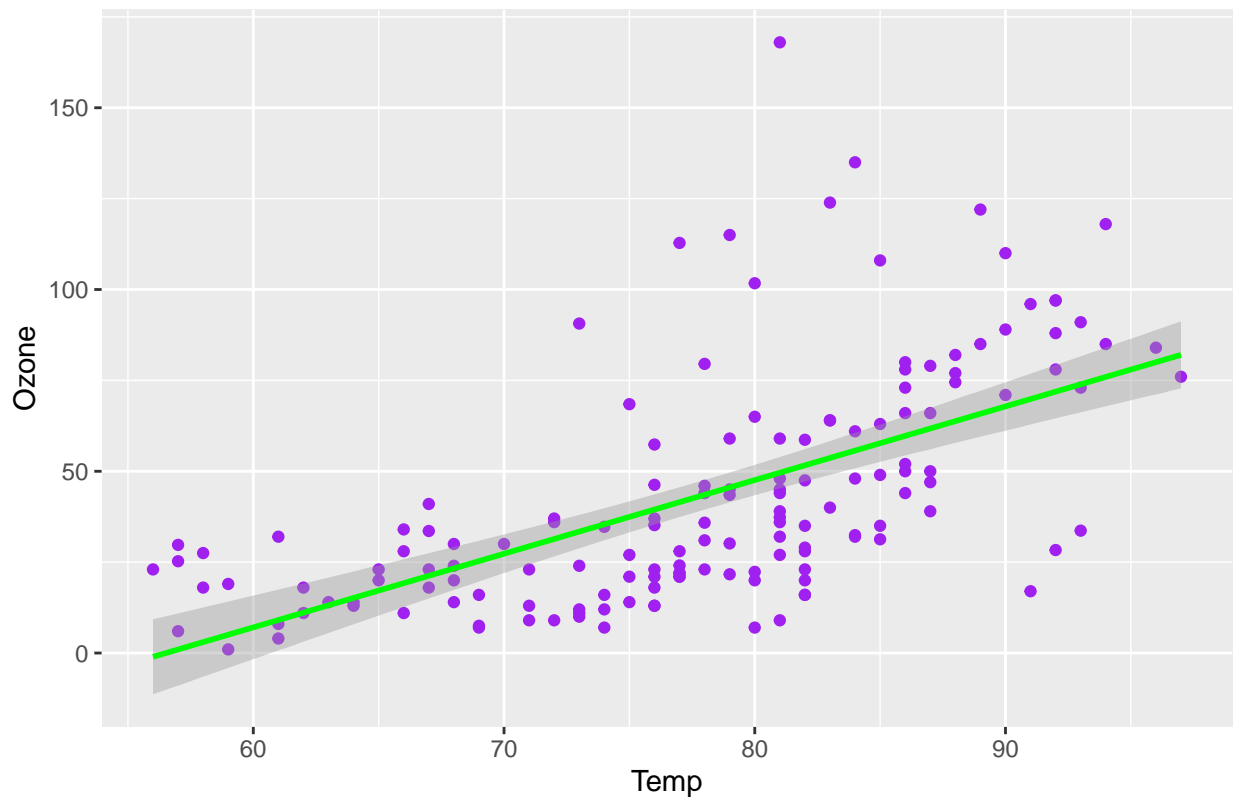
Scatterplot of Ozone by Wind



```
#Scatterplot using Ozone as outcome and Temp as predictor  
ggplot(air)+  
  geom_point(aes(Temp,Ozone),color='purple')+  
  stat_smooth(aes(Temp,Ozone),method='lm',color='green')+  
  ggtitle("Scatterplot of Ozone by Temp")
```

```
## 'geom_smooth()' using formula 'y ~ x'
```

Scatterplot of Ozone by Temp



F. Next, create a **simple regression model** predicting **Ozone** based on **Wind**, using the `lm()` command. In a comment, report the **coefficient** (aka **slope** or **beta weight**) of **Wind** in the regression output and, **if it is statistically significant, interpret it** with respect to **Ozone**. Report the **adjusted R-squared** of the model and try to explain what it means.

```
regression1 <-lm(Ozone~Wind,air)
regression1
```

```
##
## Call:
## lm(formula = Ozone ~ Wind, data = air)
##
## Coefficients:
## (Intercept)      Wind
##      89.021      -4.592
```

```
# The slope of Wind is -4.5892
```

```
summary(regression1)
```

```
##
## Call:
## lm(formula = Ozone ~ Wind, data = air)
##
```

```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -50.332 -18.332  -4.155   14.163   94.594
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   89.0205     6.6991  13.288 < 2e-16 ***
## Wind          -4.5925     0.6345  -7.238 2.15e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 27.56 on 151 degrees of freedom
## Multiple R-squared:  0.2576, Adjusted R-squared:  0.2527
## F-statistic: 52.39 on 1 and 151 DF, p-value: 2.148e-11
```

*#The adjusted R-squared is 0.2527. 25.27% of the variance in Ozone can be accounted for by Wind.
 #The p-value is 2.15e-11, which is smaller than our alpha value, 0.05. Therefore, it is statistically significant.*

- G. Create a **multiple regression model** predicting **Ozone** based on **Solar.R**, **Wind**, and **Temp**. Make sure to include all three predictors in one model – NOT three different models each with one predictor.

```
regression2 <- lm(Ozone~Solar.R+Wind+Temp,air)
summary(regression2)
```

```
##
## Call:
## lm(formula = Ozone ~ Solar.R + Wind + Temp, data = air)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -39.651 -15.622  -4.981   12.422  101.411
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -52.16596    21.90933  -2.381   0.0185 *
## Solar.R       0.01654     0.02272   0.728   0.4678
## Wind        -2.69669     0.63085  -4.275 3.40e-05 ***
## Temp         1.53072     0.24115   6.348 2.49e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 24.26 on 149 degrees of freedom
## Multiple R-squared:  0.4321, Adjusted R-squared:  0.4207
## F-statistic: 37.79 on 3 and 149 DF, p-value: < 2.2e-16
```

- H. Report the **adjusted R-Squared** in a comment – how does it compare to the adjusted R-squared from Step F? Is this better or worse? Which of the predictors are **statistically significant** in the model? In a comment, report the coefficient of each predictor that is statistically significant. Do not report the coefficients for predictors that are not significant.

```

#The adjusted R-Squared value is 0.4207. 42.07% of the variance in Ozone can be accounted for by all the
#This model is better than the one from Step F.

#Solar.R is not statistically significant
# Wind is statistically significant with a p-value of 3.40e-05
# Temp is statistically significant with a p-value of 2.49e-09

```

I. Create a one-row data frame like this:

```
predDF <- data.frame(Solar.R=290, Wind=13, Temp=61)
```

and use it with the `predict()` function to predict the **expected value of Ozone**:

```
predict(regression2, predDF, type='response')
```

```
##      1
## 10.9464
```

J. Create an additional **multiple regression model**, with **Temp** as the **outcome variable**, and the other **3 variables** as the **predictors**.

Review the quality of the model by commenting on its **adjusted R-Squared**.

```
regression3 <- lm(Temp~Ozone+Solar.R+Wind,air)
summary(regression3)
```

```
##
## Call:
## lm(formula = Temp ~ Ozone + Solar.R + Wind, data = air)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -18.831  -4.802   1.174   4.880  18.004
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  74.693222   2.796787   26.707 < 2e-16 ***
## Ozone         0.139055   0.021907    6.348 2.49e-09 ***
## Solar.R       0.015751   0.006737    2.338 0.02072 *
## Wind        -0.580176   0.195774   -2.963 0.00354 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.313 on 149 degrees of freedom
## Multiple R-squared:  0.4148, Adjusted R-squared:  0.403
## F-statistic: 35.21 on 3 and 149 DF, p-value: < 2.2e-16
```


The adjusted r -squared is 0.403. 40.3% of the variance in Temp can be accounted for by the three predi