

Linear Regression Project

Guan-Yuan Wang

2020/9/1

Linear Regression Project

For this project you will be doing the Bike Sharing Demand Kaggle challenge! We won't submit any results to the competition, but feel free to explore Kaggle more in depth. The main point of this project is to get you feeling comfortable with Exploratory Data Analysis and begin to get an understanding that sometimes certain models are not a good choice for a data set. In this case, we will discover that Linear Regression may not be the best choice given our data!

Instructions

Just complete the tasks outlined below.

Get the Data

You can download the data or just use the supplied csv in the repository. The data has the following features:

- datetime - hourly date + timestamp
- season - 1 = spring, 2 = summer, 3 = fall, 4 = winter
- holiday - whether the day is considered a holiday
- workingday - whether the day is neither a weekend nor holiday
- weather -
 - 1: Clear, Few clouds, Partly cloudy, Partly cloudy
 - 2: Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, - Mist
 - 3: Light Snow, Light Rain + Thunderstorm + Scattered clouds, - Light Rain + Scattered clouds
 - 4: Heavy Rain + Ice Pellets + Thunderstorm + Mist, Snow + Fog
- temp - temperature in Celsius
- atemp - “feels like” temperature in Celsius
- humidity - relative humidity
- windspeed - wind speed
- casual - number of non-registered user rentals initiated
- registered - number of registered user rentals initiated
- count - number of total rentals

Read in bikeshare.csv file and set it to a dataframe called bike.

```
df <- read.csv("bikeshare.csv")
```

Check the head of df

```
head(df)
```

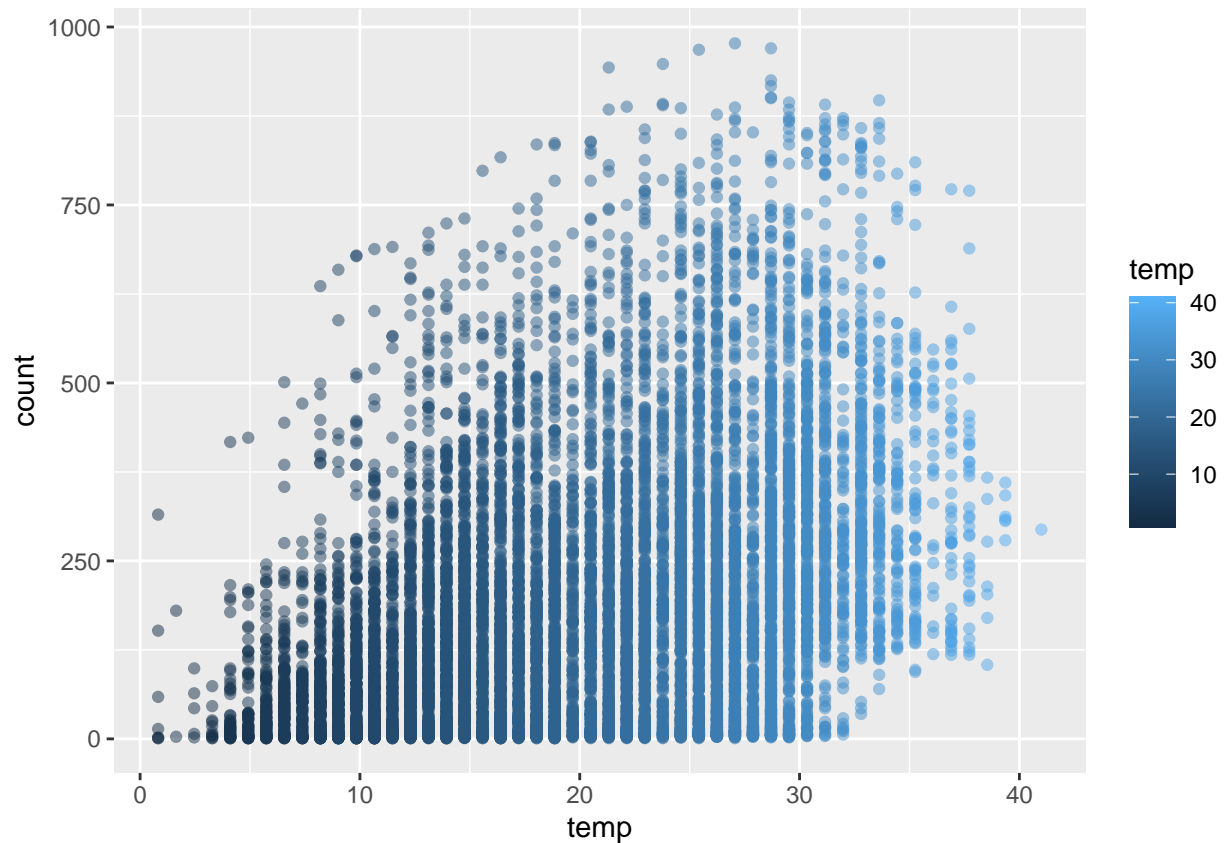
```
##           datetime season holiday workingday weather temp  atemp humidity
## 1 2011-01-01 00:00:00      1      0          0      1 9.84 14.395      81
## 2 2011-01-01 01:00:00      1      0          0      1 9.02 13.635      80
## 3 2011-01-01 02:00:00      1      0          0      1 9.02 13.635      80
## 4 2011-01-01 03:00:00      1      0          0      1 9.84 14.395      75
## 5 2011-01-01 04:00:00      1      0          0      1 9.84 14.395      75
## 6 2011-01-01 05:00:00      1      0          0      2 9.84 12.880      75
##   windspeed casual registered count
## 1    0.0000      3          13     16
## 2    0.0000      8          32     40
## 3    0.0000      5          27     32
## 4    0.0000      3          10     13
## 5    0.0000      0           1      1
## 6    6.0032      0           1      1
```

Can you figure out what is the target we are trying to predict? Check the Kaggle Link above if you are confused on this.

Exploratory Data Analysis

Create a scatter plot of count vs temp. Set a good alpha value.

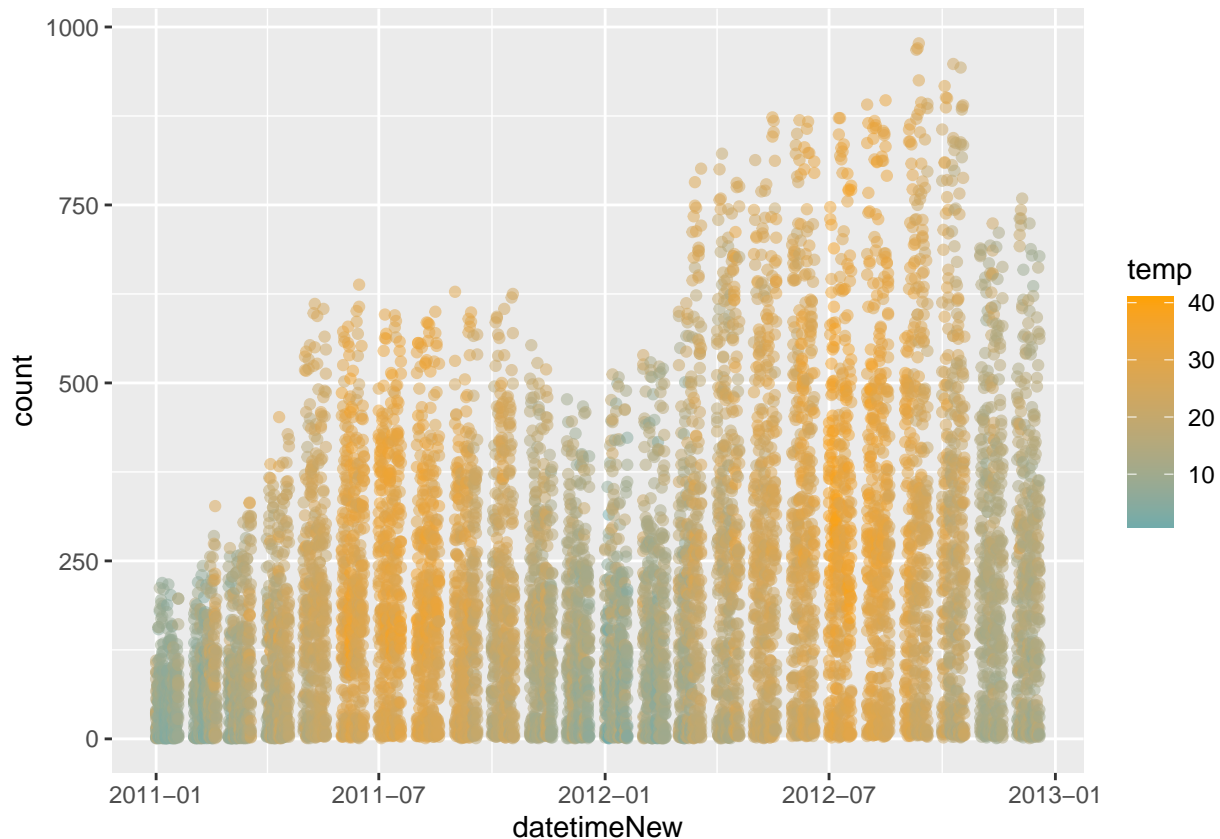
```
ggplot(df, aes(x = temp, y = count, color = temp)) +  
  geom_point(alpha = 0.5)
```



Plot count versus datetime as a scatterplot with a color gradient based on temperature. You'll need to convert the datetime column into POSIXct before plotting.

```
df$datetimeNew <- as.Date(df$datetime)

ggplot(df, aes(datetimeNew, count, color = temp)) +
  geom_point(alpha = 0.5) +
  scale_color_gradient(low = "#71abab", high = "#ffa200")
```



Hopefully you noticed two things: A seasonality to the data, for winter and summer. Also that bike rental counts are increasing in general. This may present a problem with using a linear regression model if the data is non-linear. Let's have a quick overview of pros and cons right now of Linear Regression:

Pros: - Simple to explain - Highly interpretable - Model training and prediction are fast - No tuning is required (excluding regularization) - Features don't need scaling - Can perform well with a small number of observations - Well-understood

Cons: - Assumes a linear relationship between the features and the response - Performance is (generally) not competitive with the best supervised learning methods due to high bias - Can't automatically learn feature interactions

We'll keep this in mind as we continue on. Maybe when we learn more algorithms we can come back to this with some new tools, for now we'll stick to Linear Regression.

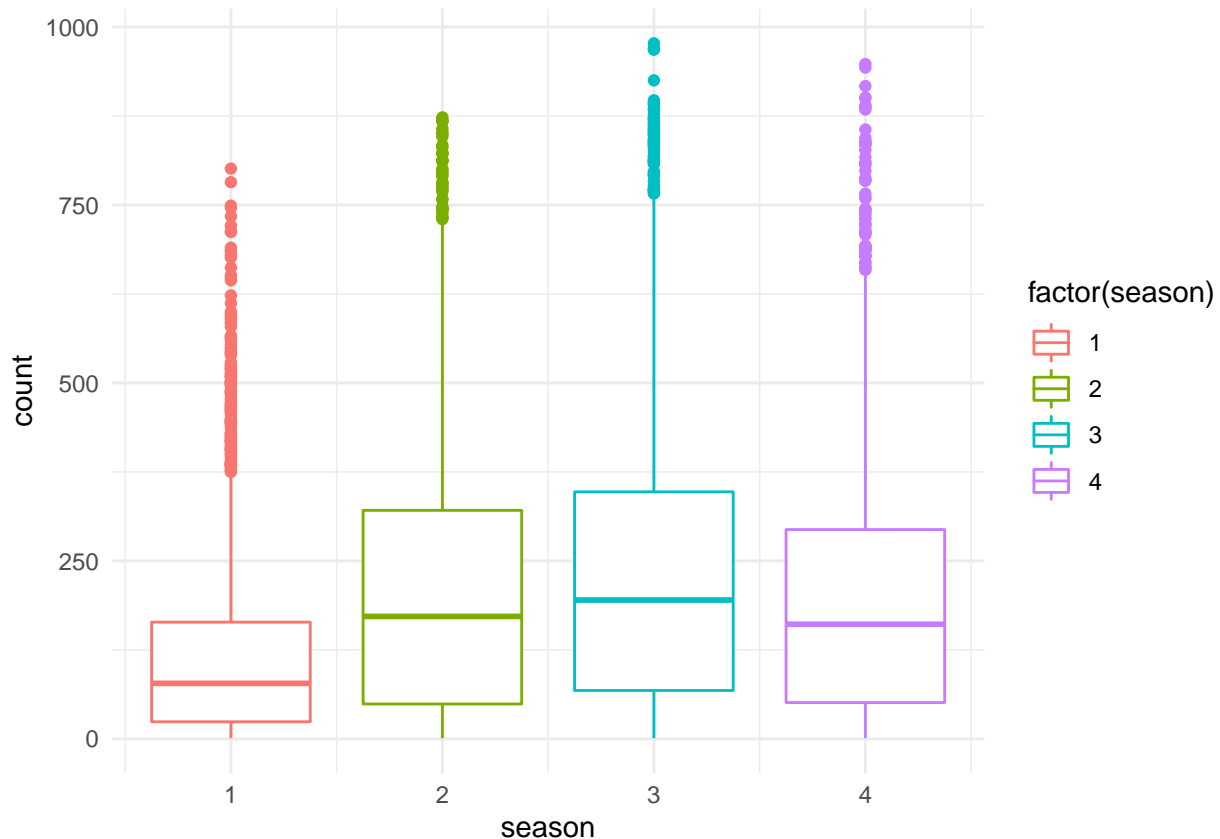
What is the correlation between temp and count?

```
cor(df$temp, df$count)
```

```
## [1] 0.3944536
```

Let's explore the season data. Create a boxplot, with the y axis indicating count and the x axis begin a box for each season.

```
ggplot(df, aes(season, count, color = factor(season))) +  
  geom_boxplot() +  
  theme_minimal()
```



Notice what this says: - A line can't capture a non-linear relationship. - There are more rentals in winter than in spring

We know of these issues because of the growth of rental count, this isn't due to the actual season!

Feature Engineering

A lot of times you'll need to use domain knowledge and experience to engineer and create new features. Let's go ahead and engineer some new features from the datetime column.

Create an "hour" column that takes the hour from the datetime column. You'll probably need to apply some function to the entire datetime column and reassign it. Hint:

```
time.stamp <- bike$datetime[4] format(time.stamp, "%H")
```

```
df$hour <- hour(df$datetime)
head(df)
```

```
##          datetime season holiday workingday weather temp  atemp humidity
## 1 2011-01-01 00:00:00     1      0         0       1  9.84 14.395      81
## 2 2011-01-01 01:00:00     1      0         0       1  9.02 13.635      80
## 3 2011-01-01 02:00:00     1      0         0       1  9.02 13.635      80
## 4 2011-01-01 03:00:00     1      0         0       1  9.84 14.395      75
## 5 2011-01-01 04:00:00     1      0         0       1  9.84 14.395      75
## 6 2011-01-01 05:00:00     1      0         0       2  9.84 12.880      75
##  windspeed casual registered count datetimeNew hour
```

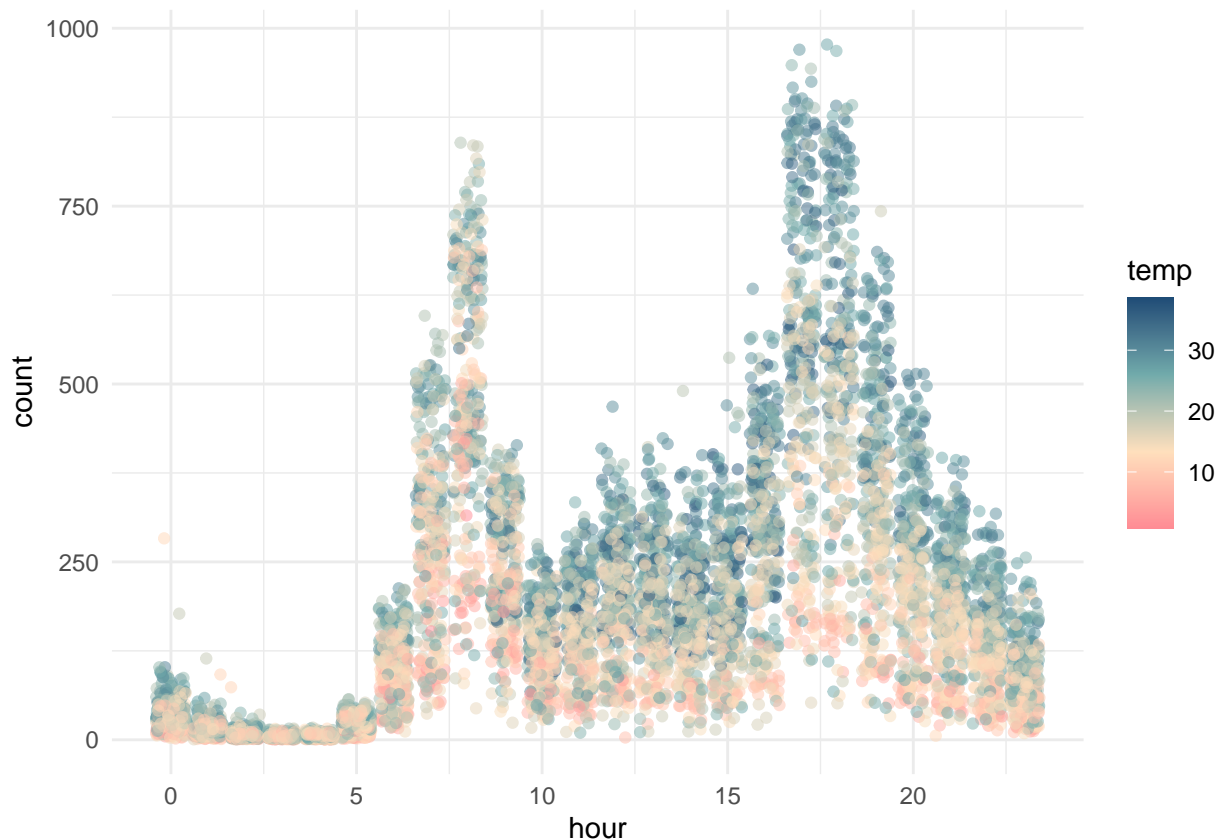
```
## 1  0.0000    3      13   16 2011-01-01    0
## 2  0.0000    8      32   40 2011-01-01    1
## 3  0.0000    5      27   32 2011-01-01    2
## 4  0.0000    3      10   13 2011-01-01    3
## 5  0.0000    0        1    1 2011-01-01    4
## 6  6.0032    0        1    1 2011-01-01    5
```

Now create a scatterplot of count versus hour, with color scale based on temp. Only use bike data where workingday==1.

Optional Additions:

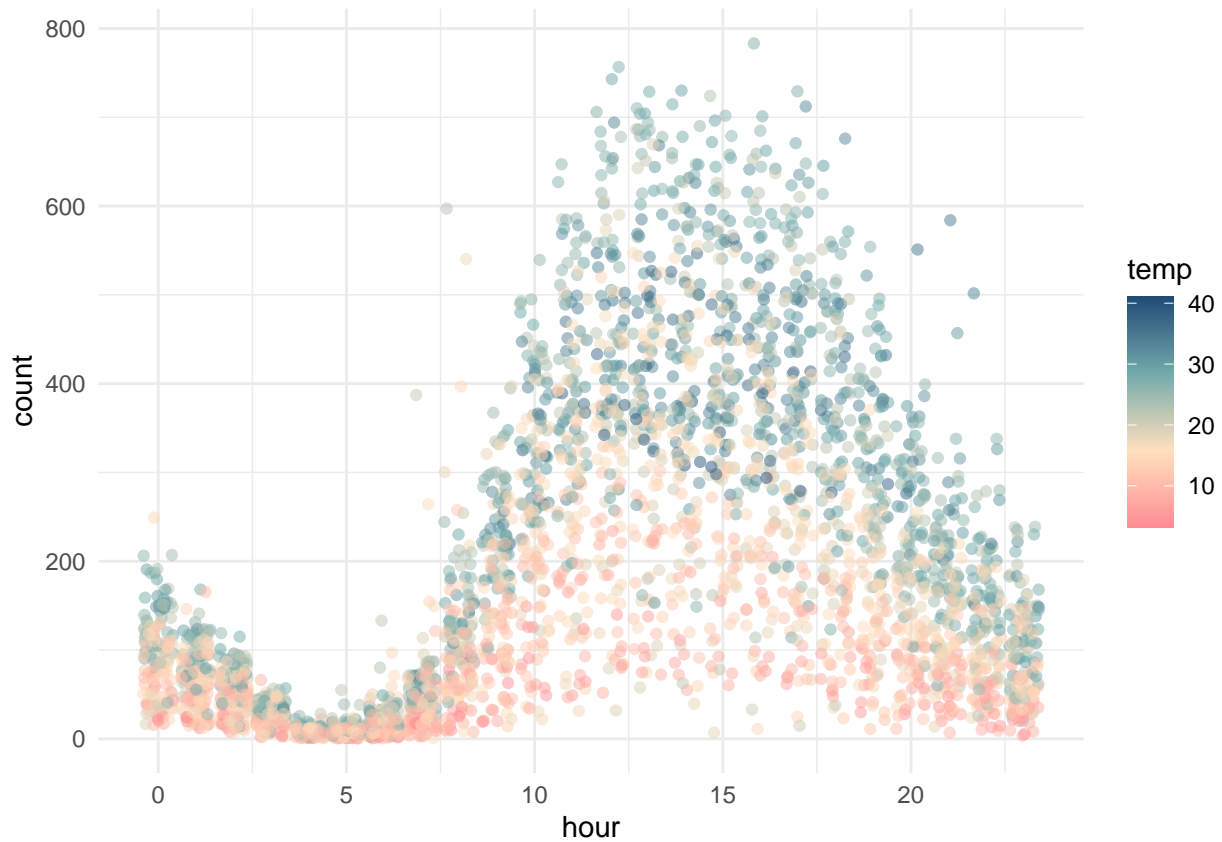
- Use the additional layer: `scale_color_gradientn(colors=c('color1',color2,etc..))` where the colors argument is a vector gradient of colors you choose, not just high and low.
- Use `position=position_jitter(w=1, h=0)` inside of `geom_point()` and check out what it does.

```
df %>%
  filter(workingday == 1) %>%
  ggplot(aes(hour, count, color = temp)) +
  geom_jitter(alpha = 0.5) +
  scale_color_gradientn(colors = c("#ff8b94", "#ffe0bd",
                                   "#71abab", "#1e4c76")) +
  theme_minimal()
```



Now create the same plot for non working days:

```
df %>%
  filter(workingday == 0) %>%
  ggplot(aes(hour, count, color = temp)) +
  geom_jitter(alpha = 0.5) +
  scale_color_gradientn(colors = c("#ff8b94", "#ffe0bd",
                                   "#71abab", "#1e4c76")) +
  theme_minimal()
```



You should have noticed that working days have peak activity during the morning (~8am) and right after work gets out (~5pm), with some lunchtime activity. While the non-work days have a steady rise and fall for the afternoon

Now let's continue by trying to build a model, we'll begin by just looking at a single feature.

Building the Model

Use `lm()` to build a model that predicts count based solely on the temp feature, name it `temp.model`

```
temp.model <- lm(count ~ temp, df)
```

Get the summary of the `temp.model`

```
summary(temp.model)
```

```
##
## Call:
## lm(formula = count ~ temp, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -293.32 -112.36  -33.36   78.98  741.44
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    6.0462     4.4394   1.362   0.173
## temp           9.1705     0.2048  44.783 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 166.5 on 10884 degrees of freedom
## Multiple R-squared:  0.1556, Adjusted R-squared:  0.1555
## F-statistic: 2006 on 1 and 10884 DF, p-value: < 2.2e-16
```

You should have gotten 6.0462 as the intercept and 9.17 as the temp coefficient. How can you interpret these values? Do some wikipedia research, re-read ISLR, or revisit the Linear Regression lecture notebook for more on this.

How many bike rentals would we predict if the temperature was 25 degrees Celsius? Calculate this two ways: - Using the values we just got above - Using the predict() function *You should get around 235.3 bikes.*

```
predict(temp.model, data.frame(temp = 25))
```

```
##      1
## 235.3097
```

Use `sapply()` and `as.numeric` to change the hour column to a column of numeric values.

```
df$hour <- sapply(df$hour, as.numeric)
```

Finally build a model that attempts to predict count based off of the following features. Figure out if theres a way to not have to pass/write all these variables into the `lm()` function. Hint: [StackOverflow](#) or [Google](#) may be quicker than the documentation. - season - holiday - workingday - weather - temp - humidity - windspeed - hour (factor)

```
model <- lm(count ~ . -casual -registered -datetime -atemp, df)
```

```
summary(model)
```

```
##
## Call:
## lm(formula = count ~ . - casual - registered - datetime - atemp,
##     data = df)
```



```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -311.72  -93.96  -28.14   62.17  646.17
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.406e+03  1.124e+02 -30.300  <2e-16 ***
## season      1.265e+00  1.458e+00   0.867   0.3858
## holiday     -9.375e+00  8.431e+00  -1.112   0.2662
## workingday  -3.451e-01  3.016e+00  -0.114   0.9089
## weather     -5.376e+00  2.396e+00  -2.244   0.0249 *
## temp        6.665e+00  1.839e-01  36.247  <2e-16 ***
## humidity    -1.997e+00  8.740e-02 -22.850  <2e-16 ***
## windspeed    4.338e-01  1.789e-01   2.425   0.0153 *
## datetimeNew  2.280e-01  7.402e-03  30.797  <2e-16 ***
## hour        7.820e+00  2.081e-01  37.574  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 141.8 on 10876 degrees of freedom
## Multiple R-squared:  0.3878, Adjusted R-squared:  0.3873
## F-statistic: 765.4 on 9 and 10876 DF,  p-value: < 2.2e-16
```

Did the model perform well on the training data? What do you think about using a Linear Model on this data?

You should have noticed that this sort of model doesn't work well given our seasonal and time series data. We need a model that can account for this type of trend, read about Regression Forests for more info if you're interested! For now, let's keep this in mind as a learning experience and move on towards classification with Logistic Regression!

Optional: See how well you can predict for future data points by creating a train/test split. But instead of a random split, your split should be "future" data for test, "previous" data for train.

Great Job!