

Collaborators:

Instructions:

Before beginning this assignment, please ensure you have access to R and/or RStudio.

1. Download the `week8a_lab.Rmd` file from Canvas. Open `week8a_lab.Rmd` in RStudio (or your favorite editor) and supply your solutions to the assignment by editing `week8a_lab.Rmd`.
2. Replace the “Insert Your Name Here” text in the `author:` field with your own full name.
3. Be sure to include code chunks, figures and written explanations as necessary. Any collaborators must be listed on the top of your assignment. Any figures should be clearly labeled and appropriately referenced within the text.
4. When you have completed the assignment and have **checked** that your code both runs in the Console and knits correctly when you click **Knit**, rename the R Markdown file to `YourLastName_YourFirstName_lab6a.Rmd`, and knit it into a PDF. Submit the compiled PDF on Canvas.

In this lab, you will need access to the following R packages:

Use the Stock Market data set to analyze the correlation of whether the market will be up or down

```
# Load some helpful libraries
library(ISLR)
```

Analyze the data

```
names(Smarket) #extract information from Smarket
```

```
## [1] "Year"      "Lag1"      "Lag2"      "Lag3"      "Lag4"      "Lag5"
## [7] "Volume"    "Today"     "Direction"
```

```
summary(Smarket) #generate results summaries
```

```
##      Year      Lag1      Lag2
## Min.   :2001  Min.   :-4.922000  Min.   :-4.922000
## 1st Qu.:2002  1st Qu.: -0.639500  1st Qu.: -0.639500
## Median :2003  Median : 0.039000  Median : 0.039000
## Mean   :2003  Mean   : 0.003834  Mean   : 0.003919
## 3rd Qu.:2004  3rd Qu.: 0.596750  3rd Qu.: 0.596750
## Max.   :2005  Max.   : 5.733000  Max.   : 5.733000
##      Lag3      Lag4      Lag5
```

```
## Min.      :-4.922000    Min.      :-4.922000    Min.      :-4.922000
## 1st Qu.: -0.640000    1st Qu.: -0.640000    1st Qu.: -0.640000
## Median : 0.038500    Median : 0.038500    Median : 0.038500
## Mean   : 0.001716    Mean   : 0.001636    Mean   : 0.00561
## 3rd Qu.: 0.596750    3rd Qu.: 0.596750    3rd Qu.: 0.59700
## Max.   : 5.733000    Max.   : 5.733000    Max.   : 5.73300
##      Volume      Today      Direction
## Min.      :0.3561    Min.      :-4.922000    Down:602
## 1st Qu.:1.2574    1st Qu.: -0.639500    Up  :648
## Median :1.4229    Median : 0.038500
## Mean   :1.4783    Mean   : 0.003138
## 3rd Qu.:1.6417    3rd Qu.: 0.596750
## Max.   :3.1525    Max.   : 5.733000
```

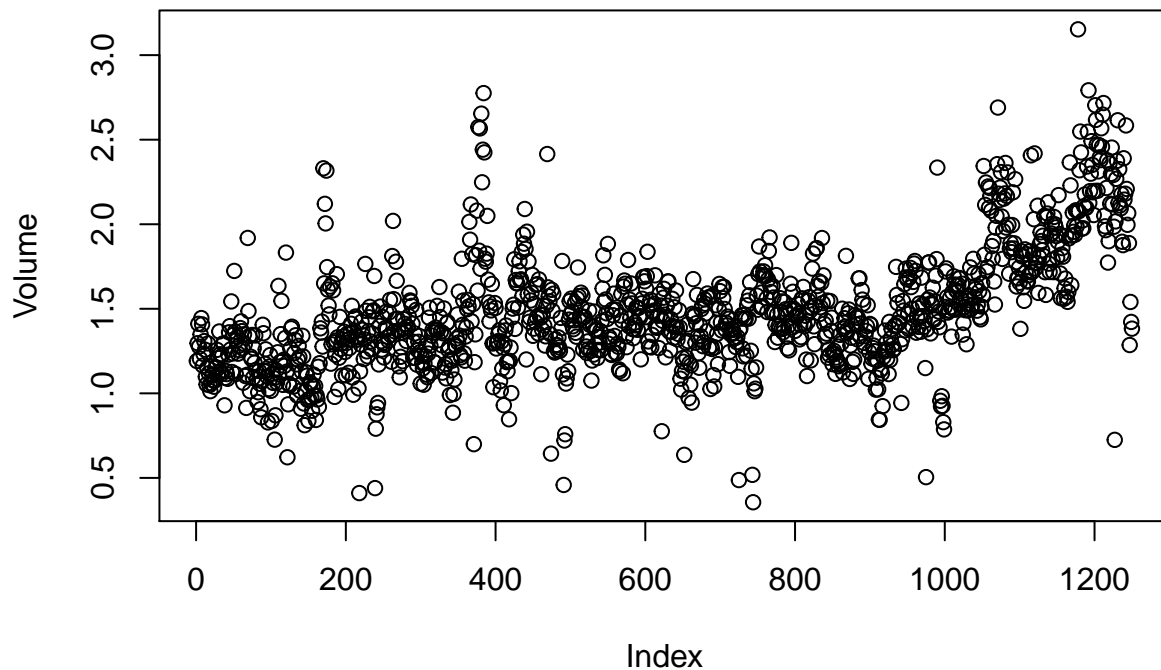
Generate a pairwise correlations

```
cor(Smarket[, -9]) #generate a matrix that contains all of the pairwise correlations
```

```
##      Year      Lag1      Lag2      Lag3      Lag4
## Year  1.00000000  0.029699649  0.030596422  0.033194581  0.035688718
## Lag1  0.02969965  1.000000000  -0.026294328  -0.010803402  -0.002985911
## Lag2  0.03059642  -0.026294328  1.000000000  -0.025896670  -0.010853533
## Lag3  0.03319458  -0.010803402  -0.025896670  1.000000000  -0.024051036
## Lag4  0.03568872  -0.002985911  -0.010853533  -0.024051036  1.000000000
## Lag5  0.02978799  -0.005674606  -0.003557949  -0.018808338  -0.027083641
## Volume 0.53900647  0.040909908  -0.043383215  -0.041823686  -0.048414246
## Today 0.03009523  -0.026155045  -0.010250033  -0.002447647  -0.006899527
##      Lag5      Volume      Today
## Year  0.029787995  0.53900647  0.030095229
## Lag1  -0.005674606  0.04090991  -0.026155045
## Lag2  -0.003557949  -0.04338321  -0.010250033
## Lag3  -0.018808338  -0.04182369  -0.002447647
## Lag4  -0.027083641  -0.04841425  -0.006899527
## Lag5  1.000000000  -0.02200231  -0.034860083
## Volume -0.022002315  1.00000000  0.014591823
## Today -0.034860083  0.01459182  1.000000000
```

Generate a graph of the volume as response

```
attach(Smarket)
plot(Volume) #plotting the correlation of volume to year
```



Logistic Regression

```
# fits generalized linear model to predict direction using lag1 through lag5 and volume
glm.fit = glm(Direction ~ Lag1+Lag2+Lag3+Lag4+Lag5+Volume, data=Smarket, family=binomial)
summary(glm.fit) #generate a summary of the generalized model
```

```
##
## Call:
## glm(formula = Direction ~ Lag1 + Lag2 + Lag3 + Lag4 + Lag5 +
##      Volume, family = binomial, data = Smarket)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.446  -1.203   1.065   1.145   1.326
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.126000   0.240736  -0.523   0.601
## Lag1        -0.073074   0.050167  -1.457   0.145
## Lag2        -0.042301   0.050086  -0.845   0.398
## Lag3         0.011085   0.049939   0.222   0.824
## Lag4         0.009359   0.049974   0.187   0.851
## Lag5         0.010313   0.049511   0.208   0.835
```

```
## Volume      0.135441  0.158360  0.855    0.392
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1731.2  on 1249  degrees of freedom
## Residual deviance: 1727.6  on 1243  degrees of freedom
## AIC: 1741.6
##
## Number of Fisher Scoring iterations: 3
```

The summary shows that the smallest p-value is associated with lag1. Thus, if we were to remove lag1 from the model then there might be a higher estimate overall.

Generate a coefficient

```
coef(glm.fit) #get the coefficient of the fitted model
```

```
## (Intercept)      Lag1      Lag2      Lag3      Lag4
## -0.12600257 -0.073073746 -0.042301344  0.011085108  0.009358938
##      Lag5      Volume
##  0.010313068  0.135440659
```

```
glm.probs = predict(glm.fit , type = "response")
```

Generate probability

```
#predicting the probability of the market going up
glm.probs = predict(glm.fit , type = "response")
glm.probs[1:10] #printing the first 10 probabilities
```

```
##      1      2      3      4      5      6      7
## 0.5070841 0.4814679 0.4811388 0.5152224 0.5107812 0.5069565 0.4926509
##      8      9     10
## 0.5092292 0.5176135 0.4888378
```

```
contrasts(Direction) #train R to generate a dummy variable for up and down... 1 = up and 0 = down
```

```
##      Up
## Down  0
## Up    1
```

```
#creating a vector of class predictions based on whether the predictability of a market increase is greater
glm.pred = rep("Down",1250)
glm.pred[glm.probs > .5] = "Up"
```

```
#determine total observations that were correctly or incorrectly classified.  
table(glm.pred , Direction)
```

```
##           Direction  
## glm.pred Down   Up  
##      Down  145 141  
##      Up   457 507
```

```
(507+145) / 1250
```

```
## [1] 0.5216
```

```
mean(glm.pred == Direction)
```

```
## [1] 0.5216
```