

# ECS 132 Final Project

Aaron Okano, Anatoly Torchinsky, Justin Maple, Samuel Huang

December 10, 2012

## 1 Forest Fire

In [Cortez and Morais, 2007], the output 'area' was first transformed with a  $\ln(x+1)$  function. Then, several Data Mining methods were applied to produce the data set of interest. Our goal with this data is to predict the fire size given a particular set of regional data.

### 1.1 Which data matters?

The dataset consists of several variables pertaining to the time, location, and weather conditions associated with each fire. There is also a variable—area—which tells us how much area the fire covered. This final variable we wish to predict using some, or all, of the remaining data.

A quick examination of the raw data reveals some important facts to us. First, there is a large quantity of observations where the area is a flat zero. Furthermore, there are very few instances of fires exceeding 100 hectares, among which several far surpass that number. This will certainly affect our regressions if left alone, so some trimming and/or splitting of the data will be necessary. Digging deeper into the data also shows some cases similar to this,

| X | Y | month | day | FFMC | DMC  | DC    | ISI  | temp | RH | wind | rain | area  |
|---|---|-------|-----|------|------|-------|------|------|----|------|------|-------|
| 3 | 4 | sep   | sun | 90.5 | 96.7 | 750.5 | 11.4 | 20.6 | 55 | 5.4  | 0    | 24.59 |
| 4 | 3 | sep   | sun | 90.5 | 96.7 | 750.5 | 11.4 | 20.4 | 55 | 4.9  | 0    | 3.64  |

which implies that there could be a high degree of variance in our estimate of the regression function.

Regardless of the above, we began our search for decent predictor variables with the raw data. This process consisted of merely running `lm` with area vs. each individual predictor variable and looking at the p-values and confidence intervals after each run. By eliminating each variable which had a coefficient estimate with a high p-value or with a confidence interval that shows little impact, the result is that not a single variable is a strong predictor. In fact, the best of the lot is temperature, which stands out with a p-value of 0.0261 and a 95% confidence interval of (0.1302, 2.0150)—a wide range consisting mostly of relatively insignificant values. Other top candidates include DMC and RH.

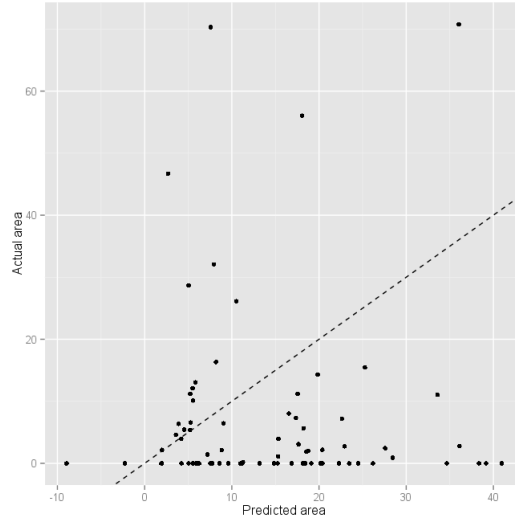
Using a forward stepping technique for adding terms to the linear model, where we continue to add terms and interaction terms to the formula and eliminate those that no longer contribute, we built the formula,

$$\hat{m}_{area;RH,DMC}(t) = -7.108 + 0.1988 \cdot t_1 + 0.316772 \cdot t_2 - 0.004806 \cdot t_1 \cdot t_2 \quad (1)$$

which is obtained from the output of `lm`.

|             | Estimate  | Std. Error | t value | Pr(> t ) |
|-------------|-----------|------------|---------|----------|
| (Intercept) | -7.107959 | 15.710450  | -0.452  | 0.6511   |
| RH          | 0.198817  | 0.312079   | 0.637   | 0.5244   |
| DMC         | 0.316772  | 0.128131   | 2.472   | 0.0138 * |
| RH:DMC      | -0.004806 | 0.002431   | -1.977  | 0.0485 * |

To test the predictive power of this particular model, we employ the method of cross validation, using 440 samples for the training set and the remaining 77 for the validation set. The result of one of these trials is the plot of actual area vs. predicted area in Figure 1.



**Figure 1:** Raw data regression—predicted area versus actual area. Dotted line is the ideal case of perfect prediction

Clearly, our current model is not particularly accurate. In fact, the  $R^2$  value for this run evaluates to a mere 0.0001555681. However, it will serve well as a benchmark for testing the predictive abilities of other models relative to this one. We now explore some possible transformations and/or restrictions on the data set in an attempt to improve our predictions.

Initially, doing variations and correlations between predictor variables and the response variable resulted in abnormally low values. This led us to consider what the forest looked like. We did this by creating a `frequencyGrid` and `areaGrid` that corresponds to each of X,Y coordinates that exist. The data shown was below:

```
> frequencyGrid(data)
      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9]
[1,]    0    0    0    0    0    0    0    0    0
[2,]   19   25    0    0    0    0    0    0    0
[3,]   10    1    1   22    0   25    2    3    0
[4,]   15   27   43   36   23    9   45    1    4
[5,]    4   20    7   25    3   49   11    4    2
[6,]    0    0    4    8    4    3    2   52    1
[7,]    0    0    0    0    0    0    0    0    0
[8,]    0    0    0    0    0    0    0    1    0
[9,]    0    0    0    0    0    0    0    0    6

> areaGrid(data)
      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]      [,8]      [,9]
[1,]      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
[2,] 11.57579 18.5060      NaN      NaN      NaN      NaN      NaN      NaN      NaN
[3,] 15.71400  0.0000  6.580000  7.858182      NaN  7.711200 13.675000  8.77000      NaN
[4,] 10.01867  5.3100  2.938372 11.039722  3.206522 16.052222 10.541556 12.18000 46.4025
[5,] 28.86750  4.6315  0.3114286 11.480400  0.000000 28.245918  7.035455  0.73250  4.0800
[6,]      NaN      NaN  0.000000 10.966250  4.405000  2.863333 43.225000 24.33269 42.8700
[7,]      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
[8,]      NaN      NaN      NaN      NaN      NaN      NaN      NaN 185.76000      NaN
[9,]      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN  0.7450
```

`frequencyGrid` and `areaGrid` reveal several trends about the forest:

1.  $\approx 44\%$  of the forest locations never had a fire

2. Locations with many fires usually have small mean areas
3. Locations with few fires usually have large mean areas

These trends suggested that certain parts of the data set were irrelevant to our prediction of fire area. In order to predict fire size in the places that actually had fires, we decided to ignore data that had particular X,Y coordinates with a mean area of NaN or frequency of 0.

## 1.2 Predictor Selection 1

For the second predictor set we decided to use the attributes that are part of the Fire Weather Index (FWI) such as FFMC, DMC, DC and ISI. Before we eliminated part of the data these values did now have a good  $R^2$  value compared to the other meteorological data such as temperature, wind, humidity and rain. However, once we eliminated the data the attributes that are part of FWI gave us a far better  $R^2$  value than the other meteorological variables. Using forward step-wise we started with FWI attributes, then we tested out different interaction terms between FWI attributes and found that  $DMC \times FMC$  and  $DC \times ISI$  greatly increased the  $R^2$  of our regression. The only non FWI attribute we decided to go with was month because we saw a strong correlation between the months and the area.

$$\begin{aligned}\hat{m}_{area;FFMC,ISI,DC,month,DMC}(t) = & -23.1151201 - 0.2744283 \cdot t_1 + 0.3185786 \cdot t_2 \\ & + 0.0122545 \cdot t_3 + 2.5926093 \cdot t_4 \\ & - 0.2495317 \cdot t_5 - 0.0006686 \cdot t_2 \cdot t_3 + 0.0027410 \cdot t_1 \cdot t_5\end{aligned}$$

which is obtained from the output of `lm`.

|             | Estimate   | Std. Error | t value | Pr(> t ) |     |
|-------------|------------|------------|---------|----------|-----|
| (Intercept) | 23.1151201 | 9.0213988  | 2.562   | 0.011419 | *   |
| FFMC        | -0.2744283 | 0.1118159  | -2.454  | 0.015301 | *   |
| ISI         | 0.3185786  | 0.1871206  | 1.703   | 0.090799 | .   |
| DC          | 0.0122545  | 0.0035439  | 3.458   | 0.000715 | *** |
| month       | 2.5926093  | 0.8779764  | 2.953   | 0.003673 | **  |
| DMC         | -0.2495317 | 0.1012981  | -2.463  | 0.014935 | *   |
| ISI:DC      | -0.0006686 | 0.0002937  | -2.277  | 0.024264 | *   |
| FFMC:DMC    | 0.0027410  | 0.0011059  | 2.479   | 0.014335 | *   |

With this regression model we got an adjusted  $R^2$  value of 0.09632 which was by far the best value we got at the time. Yet, this is still not a very good model for predicting the mean area of the fire. So then we decided to try another method of finding a good predictor set by using the backwards step-wise approach.

## 1.3 Predictor Selection 2

Backwards step-wise approach

## 1.4 Final predictor set

After testing various predictor selection methods, we came up with the following predictor set:

- FFMC
- ISI
- DMC
- Month
- DC
- DC:ISI
- DMC:FFMC

## 1.5 Cross validation

to be filled...

## 2 Parkinson's disease

The dataset was created by Max Little of the University of Oxford, in collaboration with the National Centre for Voice and Speech, Denver, Colorado, who recorded the speech signals. The main aim of the data is to discriminate healthy people from those with PD, according to "status" column which is set to 0 for healthy and 1 for PD. Our goal with this data is to try to monitor patients with Parkinson's disease remotely, by simply analyzing their voices on the phone.

### 2.1 Predictor Selection 1

We used the following:

$$m_{Y;X}(t) = P(Y = 1|X = t) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 t_1 + \dots + \beta_r t_r)}} \quad (2)$$

| predictor | class 0 | class 1 | Comments |
|-----------|---------|---------|----------|
| spread1   | 0.4     | 0.9     | Good!    |
| PPE       | 0.4     | 0.9     | Good!    |

### 2.2 Predictor Selection 2

### 2.3 Final predictor set

After testing various predictor selection methods, we came up with the following predictor set:

- spread1
- PPE

### 2.4 Cross validation

to be filled...

## 3 Beyond ECS 132

Professor Devanbu's paper, "Clones: What is that smell?" assesses the validity of the "stink" that surrounds clones. This reputation stems from the long standing belief that clone's require more project maintenance, as well as their tendency to create code bloat. One of the major problems that people face with software life cycles is maintenance costs, which can require around 80% of the total cost. As a result, Devanbu and others have invested time in research to minimize maintenance costs. One of the simplest ways to lower these costs involves reducing defects in code. Three tests were structured to analyze the relationship between clones and bugs. The first test determines the bug rate of cloned code; the second test compares the bug rate in cloned code to that of regular code; and the last test checks if prolific clone groups are buggier than the non-prolific clone groups. The results and conclusions of each test were formulated using statistical computation and analysis.

The first test was an attempt to discover to what extent cloned code contributed to bugs. The graph plots the cumulative bug convergence on the Y-axis against the clone ratio on the X-axis. According to the vertical line, which represents the average clone ratio across all snapshots, both of the projects had that about 80% of bugs have a lower clone ratio than the overall project clone ratio. The test shows that in each case the bugs contained little cloned code.

The second test finds whether clones occur more often in buggy code than elsewhere. A box and whisker plot simply illustrates for all four projects that buggy code had a lower clone ratio. This provides strong evidence that clones are not a large contributor of bugs. To further support this claim, a second table displays the adjusted  $p$  values

in which a Wilcoxon paired test is used as the null hypothesis and the alternative hypothesis is set to "snapshot clone ratio > bug clone ratio". Since the  $p$  values shown are between 0.01 and 0.05, this provides moderate evidence against the null hypothesis in favor of the alternative hypothesis. In each of the four projects, clones were not a major source of bugs. This suggests that clones are less buggy than regular code.

The third test assesses whether prolific clone groups are buggier than non-prolific clone groups. This was accomplished by finding the defect density in prolific clone groups and comparing it to those found in non-prolific clone groups. Figure 2 utilizes the resulting data in a box and whisker plot. The graph shows that in each case the bug density in prolific clone groups is lower than that of the non-prolific clone group. In addition, table three utilizes a Wilcoxon test with the alternative hypothesis set to "defect density in non-prolific groups > defect density of prolific group" provides  $p$  values that are below 0.5. These  $p$  values reject the null hypothesis in favor of the alternative, providing evidence in that more prolific clone groups are less buggy than non-prolific clone groups.

In each of the three tests, the clones were falsely attributed to a bad trait. Since the four projects used were medium to large open source projects, this data should be applicable in most situations. One area of contention over these results is how cloned code is identified. To address this issue, the tests were applied to two separate data sets. The first data set uses a conservative clone detector, while the second employs a more liberal one. Both detectors require 50 tokens in length to consider a code segment to be cloned, but the liberal detector allows for 1% less similarity than that of the conservative detector. Since analyzing both data sets led to the same conclusion, the difference in detection methods is negligible. Ultimately, the evidence from each test shows that clones may have unfairly garnered a bad reputation.

## A Code

### A.1 Problem 1

Listing 1: Analysis for Forest Fire Data Set

```
1 # Construct matrix with mean area for each coordinate
2 areaGrid <- function( data ) {
3   # For a general dataset of the same format, since I'm in the mood
4   origin <- min( data$X,data$Y )
5   maximum <- max( data$X,data$Y )
6   # Generate all possible coordinates. Use various transformations to create a
7   # list of the desired form
8   coords <- as.list( as.data.frame( t( expand.grid( origin:maximum, origin:maximum ) ) ) )
9   # Create map with mean areas
10  t( matrix( sapply( coords,
11    function(x) mean( data[ which( data$X == x[1] & data$Y == x[2] ), ]$area ) ),
12    maximum, maximum ) )
13 }
14
15 # Modification to show mean temperature per location
16 tempGrid <- function( data ) {
17   origin <- min( data$X,data$Y )
18   maximum <- max( data$X,data$Y )
19   coords <- as.list( as.data.frame( t( expand.grid( origin:maximum, origin:maximum ) ) ) )
20   # Create map with mean temperature
21   t( matrix( sapply( coords,
22     function(x) mean( data[ which( data$X == x[1] & data$Y == x[2] ), ]$temp ) ),
23     maximum, maximum ) )
24 }
25
26 # Modification of above to show number of fires per location
27 frequencyGrid <- function( data ) {
28   origin <- min( data$X,data$Y )
29   maximum <- max( data$X,data$Y )
30   coords <- as.list( as.data.frame( t( expand.grid( origin:maximum, origin:maximum ) ) ) )
31   # Create map with fire frequencies
32   t( matrix( sapply( coords,
33     function(x) length( data[ which( data$X == x[1] & data$Y == x[2] ), ]$area ) ),
34     maximum, maximum ) )
35 }
36
37 # Find the mean area for each unique entry in dataVector, e.g. the mean area
38 # for each value of DC would be found with meanArea( data, data$DC ).
39 # Returns a matrix where the first column is the value in the vector and the
40 # second is the mean area for that value.
41 meanArea <- function( data, dataVector ) {
42   uniques <- unique( dataVector )
43   cbind( uniques, lapply( uniques, function(x)
44     mean( data[ dataVector == x, ]$area ) ) )
45 }
46
47 genCond <- function( v ) {
48   sorted <- sort( v )
49   s <- split( sorted, ceiling( seq( length(v) ) / (length(v)/3) ) )
```

```

50   cbind( v >= head(s$'1', n=1) & v <= tail(s$'1', n=1),
51         v >= head(s$'2', n=1) & v <= tail(s$'2', n=1),
52         v >= head(s$'3', n=1) & v <= tail(s$'3', n=1))
53 }
54
55 # Mean area based on conditions
56 meanConds <- function(data, conditions, v) {
57   mean( data[ !apply( conditions, 1, function(x)
58                     any(!x[cbind(v,1:length(x[1,]))])), ]$area )
59 }
60
61 # Makes condition arrays using list of variables to form conditions on
62 makeConds <- function( variables ) {
63   m <- lapply( variables, genCond )
64   array( unlist( m ), dim = c(dim(m[[1]]), length(m)))
65 }
66
67 linear <- function(t,b) { b%*%c(1,t) }
68
69 cvlm <- function( variables, ntest ) {
70   v <- sample( 1:nrow(variables), ntest )
71   model <- lm( area ~ ., variables[v,] )
72   cbind( apply( as.matrix( variables[-v,-1] ), 1, linear, model$coefficients
73             ), variables[-v,1] )
74 }
75
76
77 #conds <- makeConds( data[,5:11] )
78 #perms <- data.matrix( expand.grid( rep( list(1:4), 7 ) ) )
79 #means <- apply( perms, 1, function(x) meanConds( data, conds, x ) )
80 #table <- cbind( perms, means )
81 #table <- table[complete.cases(table),]
82 #table <- table[order(table[,8]),]
83 #mapply( function(x,y,n) mean( table[ table[,8] > x & table[,8] <= y, ][,n] ),
84 #       0, 4.4, 1:7 )
85
86
87 #model <- lm( log( area + 1 ) ~ month + temp + RH + DC + month:RH, data=ordata)
88
89
90 #predMat <- cbind( ordata$X, ordata$month, ordata$FFMC, ordata$DC, ordata$X *
91 #               ordata$FFMC )
92
93 #mean( sapply( 1:4000, function(n) {v <- sample(1:153, 150); cor(
94 #               newdata$area[-v],
95 #               knn(
96 #                 predMat[v,],
97 #                 newdata$area[v],
98 #                 10,
99 #                 predMat[-v,]
100 #               )$predyvals
101 #             )^2 )
102
103

```

```

104 #> summary( lm( area ~ FPMC + ISI + ISI:FPMC + DC + DC:ISI + month, data=ordata
105 #           > ) )
106 #
107 #Call:
108 #lm(formula = area ~ FPMC + ISI + ISI:FPMC + DC + DC:ISI + month,
109 #    data = ordata)
110 #
111 #Residuals:
112 #      Min       1Q   Median       3Q      Max
113 #-3.353 -1.556 -0.476  1.118  4.743
114 #
115 #Coefficients:
116 #              Estimate Std. Error t value Pr(>|t|)
117 #(Intercept)  18.7922107   8.5250001   2.204  0.02909 *
118 #FFMC         -0.2123159   0.1031263  -2.059  0.04132 *
119 #ISI          -2.1986055   1.3395298  -1.641  0.10291
120 #DC            0.0103264   0.0031436   3.285  0.00128 **
121 #month         2.6778341   0.8689659   3.082  0.00247 **
122 #FFMC:ISI      0.0258508   0.0146401   1.766  0.07956 .
123 #ISI:DC       -0.0004180   0.0002632  -1.588  0.11442
124 #---
125 #Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
126 #
127 #Residual standard error: 2.007 on 144 degrees of freedom
128 #Multiple R-squared:  0.1142, Adjusted R-squared:  0.0773
129 #F-statistic: 3.094 on 6 and 144 DF, p-value: 0.007022
130
131 ##### Better Result #####
132
133 #> g <- lm(formula = area ~ FPMC + ISI + DC + DC:ISI + month + DMC + DMC:FFMC , data=ordata,
134 #> summary(g)
135
136 #Call:
137 #lm(formula = area ~ FPMC + ISI + DC + DC:ISI + month + DMC +
138 #    DMC:FFMC, data = ordata)
139
140 #Residuals:
141 #      Min       1Q   Median       3Q      Max
142 #-3.3620 -1.4821 -0.3396  1.0433  4.8055
143
144 #Coefficients:
145 #              Estimate Std. Error t value Pr(>|t|)
146 #(Intercept)  23.1151201   9.0213988   2.562  0.011419 *
147 #FFMC         -0.2744283   0.1118159  -2.454  0.015301 *
148 #ISI          0.3185786   0.1871206   1.703  0.090799 .
149 #DC            0.0122545   0.0035439   3.458  0.000715 ***
150 #month         2.5926093   0.8779764   2.953  0.003673 **
151 #DMC          -0.2495317   0.1012981  -2.463  0.014935 *
152 #ISI:DC       -0.0006686   0.0002937  -2.277  0.024264 *
153 #FFMC:DMC      0.0027410   0.0011059   2.479  0.014335 *
154 #---
155 #Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
156
157 #Residual standard error: 1.999 on 145 degrees of freedom

```



```

158 #Multiple R-squared: 0.1379, Adjusted R-squared: 0.09632
159 #F-statistic: 3.314 on 7 and 145 DF, p-value: 0.002658
160
161
162 #> predictorMat <- cbind( ordata$FFMC, ordata$ISI, ordata$DC, ordata$month, ordata$DMC, ordata$
163 data <- read.csv( 'forestfires.csv', head=TRUE)
164 data$month <- factor( data$month,
165     levels=c( 'jan', 'feb', 'mar', 'apr', 'may', 'jun',
166         'jul', 'aug', 'sep', 'oct', 'nov', 'dec' ))
167 data$day <-
168     factor( data$day, levels=c( 'mon', 'tue', 'wed', 'thu', 'fri', 'sat', 'sun' ))
169 #data$month <- sin( as.integer( data$month ) * pi / 6 )
170 #data$day <- sin( as.integer( data$day ) * 2 * pi / 7 )
171 ordata <- data[ order( data$area ), ][248:400, ]
172 ordata2 <- data[ order( data$area ), ][248:448, ]
173 ordata2$area <- log( ordata2$area + 1 )

```

## A.2 Problem 2

**Listing 2:** Analysis for Parkinson's Data Set

```

1 parkinson <- read.csv( 'parkinsons.data', header=TRUE)
2
3 # Generalized logit function. Inputs are the input variables (the t's in the
4 # book) and the coefficients (the betas in the book). Both inputs are in vector
5 # form. Naturally, the vector t should have one less element than the vector b.
6 logit <- function( t, b ) { 1/(1+exp(-(b %*% c(1,t)))) }
7 crossvalglm <- function( response, predictor, predictor2, predictor3 ) {
8     v <- sample( 1:(length(response) - 1), (length( response ) - 1) * 0.5 )
9     notv <- setdiff( 1:(length(response) - 1), v )
10    model <- glm( response[v] ~ predictor[v] + predictor2[v] + predictor3[v],
11        family = binomial )
12    cor( response[notv], mapply( logit3, model$coefficients[1],
13        model$coefficients[2],
14        model$coefficients[3],
15        model$coefficients[4],
16        predictor[notv],
17        predictor2[notv],
18        predictor3[notv])
19    )^2
20 }
21
22 crossvalglm2 <- function( response, predictor, predictor2 ) {
23     v <- sample( 1:(length(response) - 1), (length( response ) - 1) * 0.5 )
24     notv <- setdiff( 1:(length(response) - 1), v )
25     model <- glm( response[v] ~ predictor[v] + predictor2[v], family = binomial
26         )
27     cor( response[notv], mapply( logit2, model$coefficients[1],
28        model$coefficients[2],
29        model$coefficients[3],
30        predictor[notv],
31        predictor2[notv])
32    )^2

```

```

33 }
34
35 crossvalglm3 <- function( response, predictor ) {
36   v <- sample( 1:(length(response) - 1), (length( response ) - 1) * 0.5 )
37   notv <- setdiff( 1:(length(response) - 1), v )
38   model <- glm( response[v] ~ predictor[v,], family = binomial )
39   cor( response[notv], mapply( logit , model$coefficients[1],
40                               model$coefficients[2],
41                               predictor[notv]
42                               ) )^2
43 }
44
45 # Generalized glm cross validation. Takes as input a data frame where the first
46 # column is status and the remaining columns are predictors. Returns proportion
47 # of successful predictions.
48 # Call the function with cvglmprop( parkinson[,c("status","var1","var2",...)] )
49 # If you want to add interaction terms, perform the multiplication beforehand:
50 # temp <- parkinson[,c("status","var1","var2")]
51 # temp$var1v2 <- temp$var1 * temp$var2
52 # cvglmprop( temp )
53 cvglmprop <- function( variables ) {
54   v <- sample( 1:(nrow(variables) - 1), (nrow( variables ) - 1) * 0.5 )
55   notv <- setdiff( 1:(nrow(variables) - 1), v )
56   model <- glm( status ~ ., data = variables[v,], family = binomial )
57   pred <- apply( as.matrix( variables[notv,-1] ), 1, logit ,
58                 model$coefficients )
59   pred <- pred > 0.5
60   mean( pred == variables[notv,1] )
61 }

```

## B Who did what

- Aaron - implemented R functions for Problem 1 and 2 and did initial analysis
- Anatoly - assisted Aaron in development of Problem 1 and 2 code; used code to generate plots
- Justin - used Devanbu's article to finish Problem 3
- Samuel - wrote R to do initial analysis, copy-edited Problem 3 and did write-up for Problem 1 and 2