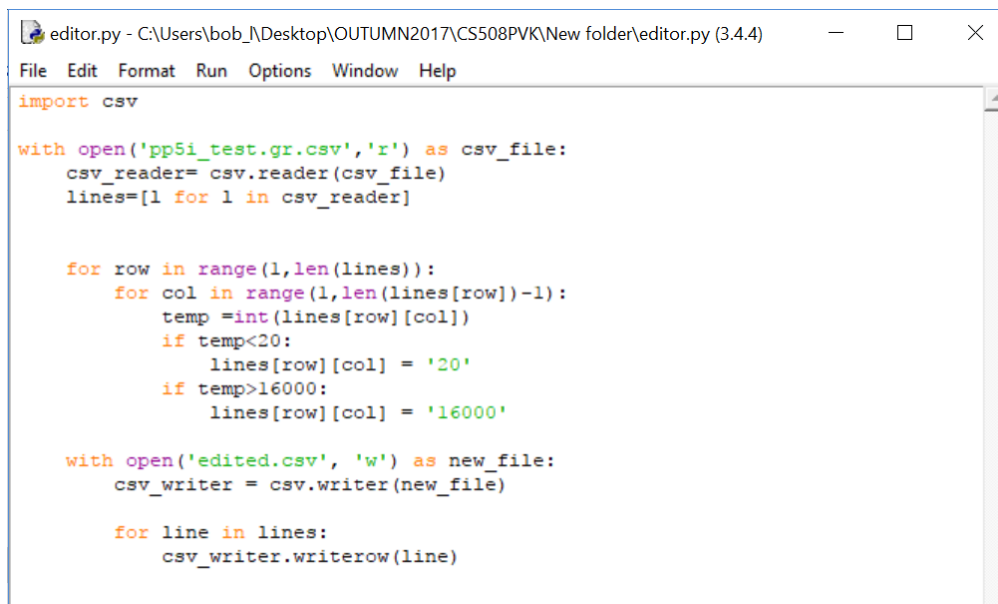**Sean Reeves**

**Jose Jarquin**

**Jeremy Keys**

**(1)**

We completed section 1 by re-using the normalize-data script we each had to write for HW5.



```python
import csv

with open('pp5i_test.gr.csv','r') as csv_file:
    csv_reader= csv.reader(csv_file)
    lines=[l for l in csv_reader]


    for row in range(1,len(lines)):
        for col in range(1,len(lines[row])-1):
            temp =int(lines[row][col])
            if temp<20:
                lines[row][col] = '20'
            if temp>16000:
                lines[row][col] = '16000'

    with open('edited.csv', 'w') as new_file:
        csv_writer = csv.writer(new_file)

        for line in lines:
            csv_writer.writerow(line)
```

*Figure 1 A simple normalize script*

**(2)**

For the first section of step (2), we re-purposed the fold calculating scripts we had to write. (For this part, we used Sean's script, which is written in R.)

```r
value <- read.csv(file =
"C:/Users/Sreav/Desktop/FinalData/pp5i_train.gr.normalized.removeless2.csv", head = TRUE,
sep=",")

str(value)
sumsq <-
 function(X){
   # sum of squares of input object

   if(is.list(X)){
     xss <- sum(sapply(X,sumsq))
   } else {
     xss <- sum(X^2)
   }

   xss
```

```
  }

for(i in 1:nrow(value)){
  avMED <- rowMeans(value[i,2:40])
  sumMED <- rowSums(value[i,2:40])
  sumSQMED <- sumsq(value[i,2:40])
  stdevMED <- sqrt((39*sumSQMED - sumMED*sumMED)/(39*(39-1)))

  avMGL <- rowMeans(value[i,41:47])
  sumMGL <- rowSums(value[i,41:47])
  sumSQMGL <- sumsq(value[i,41:47])
  stdevMGL <- sqrt((7*sumSQMGL - sumMGL *sumMGL)/(7*(7-1)))

  avRHB <- rowMeans(value[i,48:54])
  sumRHB <- rowSums(value[i,48:54])
  sumSQRHB <- sumsq(value[i,48:54])
  stdevRHB <- sqrt((7*sumSQRHB - sumRHB *sumRHB )/(7*(7-1)))

  avEPD <- rowMeans(value[i,55:64])
  sumEPD <- rowSums(value[i,55:64])
  sumSQEPD <- sumsq(value[i,55:64])
  stdevEPD <- sqrt((10*sumSQEPD - sumEPD *sumEPD )/(10*(10-1)))

  avJPA <- rowMeans(value[i,65:70])
  sumJPA <- rowSums(value[i,65:70])
  sumSQJPA <- sumsq(value[i,65:70])
  stdevJPA <- sqrt((6*sumSQJPA - sumJPA *sumJPA )/(6*(6-1)))

  avOtherMED <- (avMGL + avRHB + avEPD + avJPA)/4
  avOtherMGL <- (avMED + avRHB + avEPD + avJPA)/4
  avOtherRHB <- (avMED + avMGL + avEPD + avJPA)/4
  avOtherEPD <- (avMED + avMGL + avRHB + avJPA)/4
  avOtherJPA <- (avMED + avMGL + avRHB + avEPD)/4

  sumOtherMED <- sumMGL + sumRHB + sumEPD + sumJPA
  sumOtherMGL <- sumMED + sumRHB + sumEPD + sumJPA
  sumOtherRHB <- sumMGL + sumMED + sumEPD + sumJPA
  sumOtherEPD <- sumMGL + sumMED + sumRHB + sumJPA
  sumOtherJPA <- sumMGL + sumMED + sumRHB + sumEPD

  sumSQOtherMED <- sumSQMGL + sumSQRHB + sumSQEPD + sumSQJPA
  sumSQOtherMGL <- sumSQMED + sumSQRHB + sumSQEPD + sumSQJPA
  sumSQOtherRHB <- sumSQMGL + sumSQMED + sumSQEPD + sumSQJPA
  sumSQOtherEPD <- sumSQMGL + sumSQMED + sumSQRHB + sumSQJPA
  sumSQOtherJPA <- sumSQMGL + sumSQMED + sumSQRHB + sumSQEPD

  stdevOtherMED <- sqrt((30*sumSQOtherMED - sumOtherMED * sumOtherMED)/(30*(30-1)))
```

```r
stdevOtherMGL <- sqrt((62*sumSQOtherMGL - sumOtherMGL * sumOtherMGL)/(62*(62-1)))
stdevOtherRHB <- sqrt((62*sumSQOtherRHB - sumOtherRHB * sumOtherRHB)/(62*(62-1)))
stdevOtherEPD <- sqrt((59*sumSQOtherEPD - sumOtherEPD * sumOtherEPD)/(59*(59-1)))
stdevOtherJPA <- sqrt((63*sumSQOtherJPA - sumOtherJPA * sumOtherJPA)/(63*(63-1)))

#TvalueALL <- (avALL - avAML)/sqrt(stdevALL*stdevALL / 27 + stdevAML*stdevAML / 11)
#TvalueAML <- (avAML - avALL)/sqrt(stdevAML*stdevAML / 11 + stdevALL*stdevALL / 27)
TvalueMED <- (avMED - avOtherMED)/sqrt(stdevMED*stdevMED / 39 +
stdevOtherMED*stdevOtherMED / 30)
TvalueMED <- abs(TvalueMED)
TvalueMGL <- (avMGL - avOtherMGL)/sqrt(stdevMGL*stdevMGL / 7 +
stdevOtherMGL*stdevOtherMGL / 62)
TvalueMGL <- abs(TvalueMGL)
TvalueRHB <- (avRHB - avOtherRHB)/sqrt(stdevRHB*stdevRHB / 7 +
stdevOtherRHB*stdevOtherRHB / 62)
TvalueRHB <- abs(TvalueRHB)
TvalueEPD <- (avEPD - avOtherEPD)/sqrt(stdevEPD*stdevEPD / 10 +
stdevOtherEPD*stdevOtherEPD / 59)
TvalueEPD <- abs(TvalueEPD)
TvalueJPA <- (avJPA - avOtherJPA)/sqrt(stdevJPA*stdevJPA / 6 +
stdevOtherJPA*stdevOtherJPA / 63)
TvalueJPA <- abs(TvalueJPA)

write(TvalueMED,file="C:/Users/Sreav/Desktop/FinalData/TvalueMED.txt",append=TRUE)
write(TvalueMGL,file="C:/Users/Sreav/Desktop/FinalData/TvalueMGL.txt",append=TRUE)
write(TvalueRHB,file="C:/Users/Sreav/Desktop/FinalData/TvalueRHB.txt",append=TRUE)
write(TvalueEPD,file="C:/Users/Sreav/Desktop/FinalData/TvalueEPD.txt",append=TRUE)
write(TvalueJPA,file="C:/Users/Sreav/Desktop/FinalData/TvalueJPA.txt",append=TRUE)


}
```

After generating the top 30 genes for each class sorted in descending order, we generated the subsets by using the `head` command on Linux to retain the first `N` lines (genes) for each subset (2,4,6,8,10,12,15,20,25) and each class, `cat` to merge the resulting subset files for each N, and 'gawk' to keep only the unique lines for each subset in their original order. E.g.:

```
head -25 Top30MED.csv > Top25MED.csv
cat Top25MED Top25JPA.csv Top25MGL.csv Top25RHB.csv Top25EPD.csv > Top25-not-
unique.csv
gawk '!seen[$0]++' Top25-not-unique.csv
cat Top25-not-unique.csv | awk '!seen[$0]++' > Top25-unique.csv
```

Finally, we transposed the class file we were given, and manually added it to the last line of each top-N-unique files, then transposed each file again to get it in genes-in-columns format.

```python
#https://unix.stackexchange.com/questions/60590/is-there-a-command-
line-utility-to-transpose-a-csv-file

import csv, sys

#constants
DEBUG = True
INPUT_FILE_NAME_TRAIN_30 = "pp5i_train.top30.gr.unique.csv"
OUTPUT_FILE_NAME_TRAIN_30 = "pp5i_train.top30.gr.unique.transposed.csv"

def transposeRows(input_file_name, output_file_name):
    with open(input_file_name) as file:
        out_file = open(output_file_name, "w")
        rows = list(csv.reader(file))
        writer = csv.writer(out_file)
        for col in range(0, len(rows[0])):
            writer.writerow([row[col] for row in rows])

transposeRows(INPUT_FILE_NAME_TRAIN_30, OUTPUT_FILE_NAME_TRAIN_30)
```

**(3)**

For part 3 we selected the classifier ZeroR at the beginning, but after measuring the classifier's accuracy, we found out that with this classifier, we would get the same value for all the previously generated files and we decided to change it to a more accurate classifier.

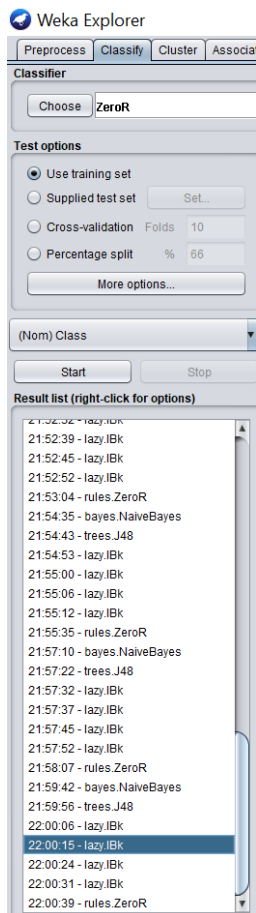The one we ended up selecting is the Multilayer Perceptron classifier.



*Figure 2 All the models we got from calculating the accuracy*

3b)

After this we calculated the cross-validation error rate the same way as before, but this time instead of using the training set, we used cross-validation with 10 folds. (See Figure 2.)
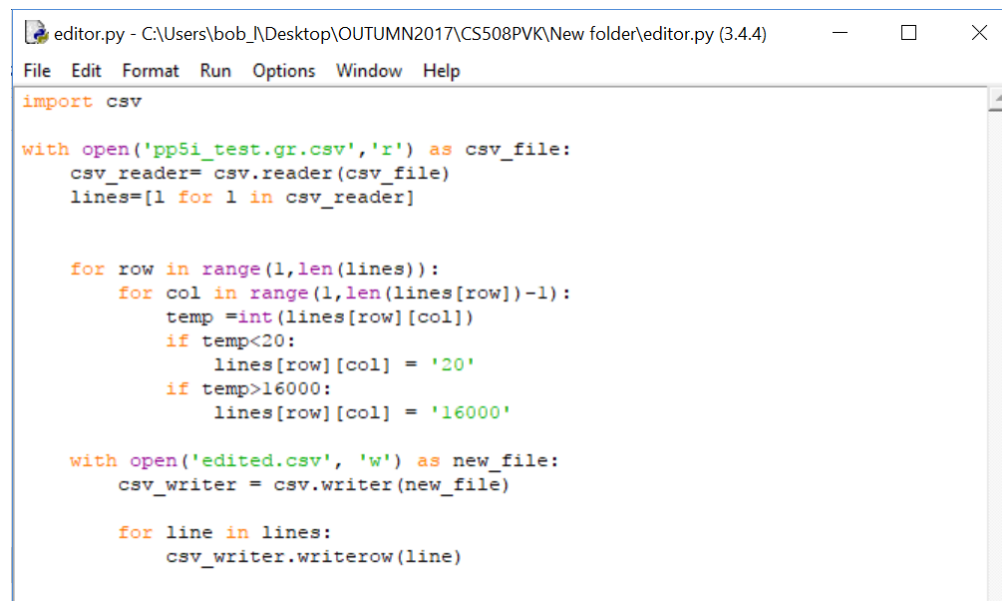
| classifier | | top2 | top4 | top6 | top8 | top10 | top12 | top15 | top20 | top25 | top30 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Naïve Bayes | accuracy | 65 | 66 | 67 | 67 | 67 | 66 | 66 | 67 | 67 | 67 |
| | error | 4 | 3 | 2 | 2 | 2 | 3 | 3 | 2 | 2 | 2 |
| J48 | accuracy | 57 | 55 | 54 | 56 | 62 | 63 | 64 | 64 | 66 | 67 |
| | error | 12 | 14 | 15 | 13 | 7 | 6 | 5 | 5 | 3 | 2 |
| IB1 | accuracy | 64 | 66 | 68 | 65 | 66 | 67 | 67 | 68 | 67 | 66 |
| | error | 5 | 3 | 1 | 1 | 3 | 2 | 2 | 1 | 2 | 3 |
| IB2 | accuracy | 62 | 63 | 65 | 64 | 61 | 64 | 65 | 66 | 66 | 66 |
| | error | 7 | 6 | 4 | 5 | 8 | 5 | 4 | 3 | 3 | 3 |
| IB3 | accuracy | 65 | 66 | 67 | 68 | 67 | 66 | 66 | 68 | 67 | 67 |
| | error | 4 | 3 | 2 | 1 | 2 | 3 | 3 | 1 | 2 | 2 |
| IB4 | accuracy | 63 | 65 | 67 | 67 | 67 | 66 | 66 | 68 | 66 | 65 |
| | error | 6 | 4 | 2 | 2 | 2 | 3 | 3 | 1 | 3 | 4 |
| Multilayer Perceptrion | accuracy | 65 | 67 | 68 | 67 | 67 | 68 | 68 | 68 | 68 | 68 |
| | error | 4 | 2 | 1 | 2 | 2 | 1 | 1 | 1 | 1 | 1 |

*Figure 3 Results obtained by running 10-fold cross-validation on each classifier.*

With this information, we decided we would select Multilayer Perceptrion because this it tied for the most accurate classifier for all of the files (6 out of 10 have error rate 1/69) and gene set of top 20 for the same reason (4 out of 7 have error rate 1/69).

3c)

Then, with the test set, we used a python program that converted all the values smaller than 20 to 20 and the values larger that 1600 to 1600.

```
editor.py - C:\Users\bob_l\Desktop\OUTUMN2017\CS508PVK\New folder\editor.py (3.4.4)        —    □    ×

File  Edit  Format  Run  Options  Window  Help

import csv

with open('pp5i_test.gr.csv','r') as csv_file:
    csv_reader= csv.reader(csv_file)
    lines=[l for l in csv_reader]


    for row in range(1,len(lines)):
        for col in range(1,len(lines[row])-1):
            temp =int(lines[row][col])
            if temp<20:
                lines[row][col] = '20'
            if temp>16000:
                lines[row][col] = '16000'

    with open('edited.csv', 'w') as new_file:
        csv_writer = csv.writer(new_file)

        for line in lines:
            csv_writer.writerow(line)
```

*Figure 4 A simpler normalize script*

3d)

To get the "top 20 train set" names so we could extract the data from the test set, we selected the "top 20 train set" data and transposed it. After that, we cut the names of the top 20 and pasted them next to the original test set file (see Figure 6).
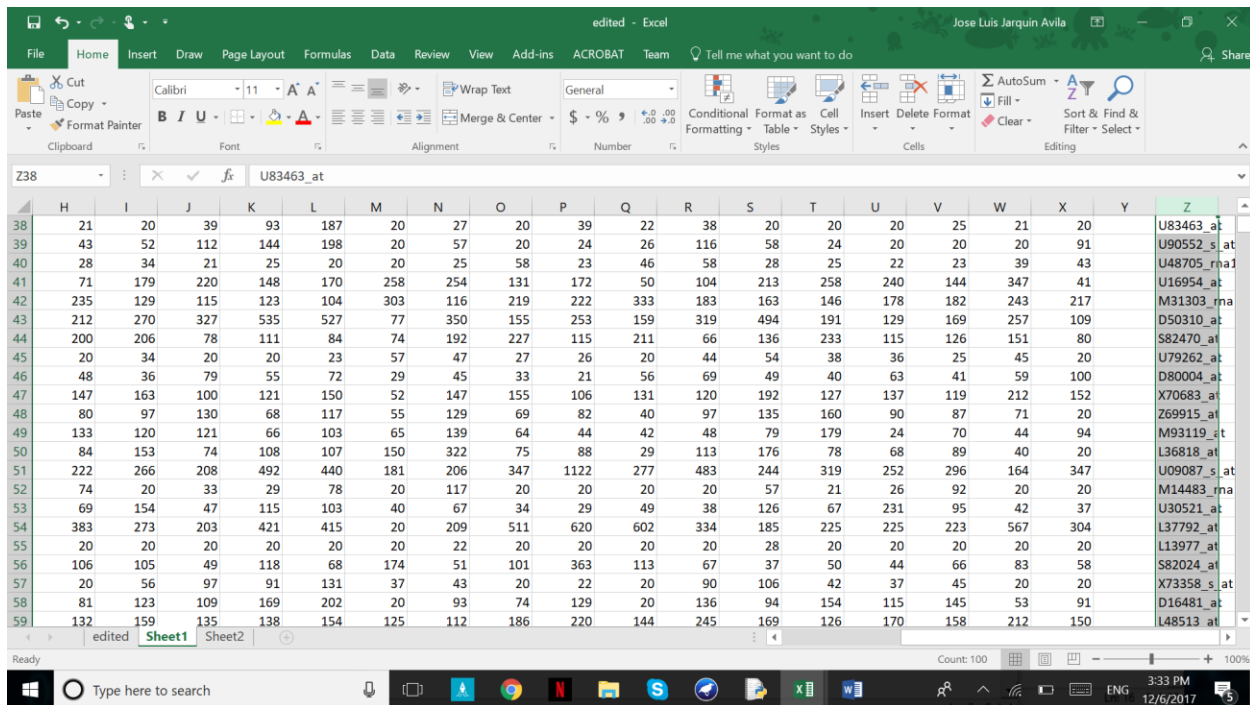
*Figure 5 Selecting unique genes*

We used an Excel formula to compare two arbitrary columns and highlight the names they have in common (the intersection of the two). First, we selected the A column as our target and the recently pasted names as out search and used the formula

**=countif($Z:$Z,$A1)**

$Z:$Z represents the entire Z (gene name) column, which is where our train set data is located, and $A represents the row where we will search. (A1 tells the function to begin at row 1.) The reason why 1 does not have a $ before it is that without it, the formula will continue to A2 and so on.



*Figure 6 Removing extraneous genes with Excel*

After the similar values were selected, we removed the ones that were not underlined. Finally, we transposed the data to have the names in columns by copying it and using "paste special" in Excel.

| | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | D13631_s | D16481_a | D25304_a | D42041_a | D50310_a | D80004_a | D83646_a | D85815_a | D87465_a | D87470_a | HG627-HT | HG880 |
| 2 | 44 | 22 | 181 | 187 | 542 | 148 | 74 | 342 | 166 | 155 | 320 | |
| 3 | 165 | 178 | 397 | 124 | 592 | 147 | 39 | 182 | 76 | 89 | 201 | |
| 4 | 20 | 20 | 20 | 76 | 953 | 77 | 62 | 108 | 700 | 77 | 274 | |

*Figure 7 transposing the genes to genes-in-columns format*

3e) We manually added the class with "?" values in Excel.

| | CU | CV |
|---|---|---|
| | _at Z80788_at | Class |
| 81 | 42 | ? |
| 17 | 28 | ? |
| 92 | 20 | ? |
| 97 | 38 | ? |
| 80 | 26 | ? |
| 84 | 61 | ? |
| 75 | 40 | ? |
| 36 | 20 | ? |
| 47 | 38 | ? |
| 61 | 33 | ? |
| 81 | 34 | ? |
| 59 | 53 | ? |
| 26 | 20 | ? |

*Figure 8 Inserting the ? column*

**(4) Generate Predictions For the Test Set**

a) The test file was converted to ARFF format by opening it in weka and then saving it.

B0We used Word to change the class line in test file to match the train file.
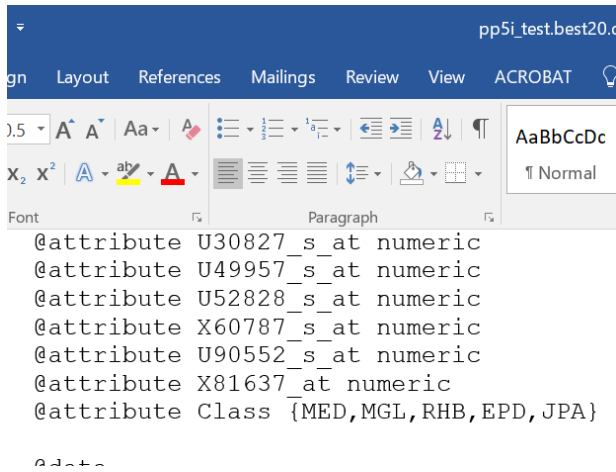


*Figure 9 Adding the matching attribute line from the train to the test file*

Per the assignment instructions, we used Weka to generate predictions by using its Preprocess>Classify tab. We added the predictions with Excel, and the *?* class by hand.
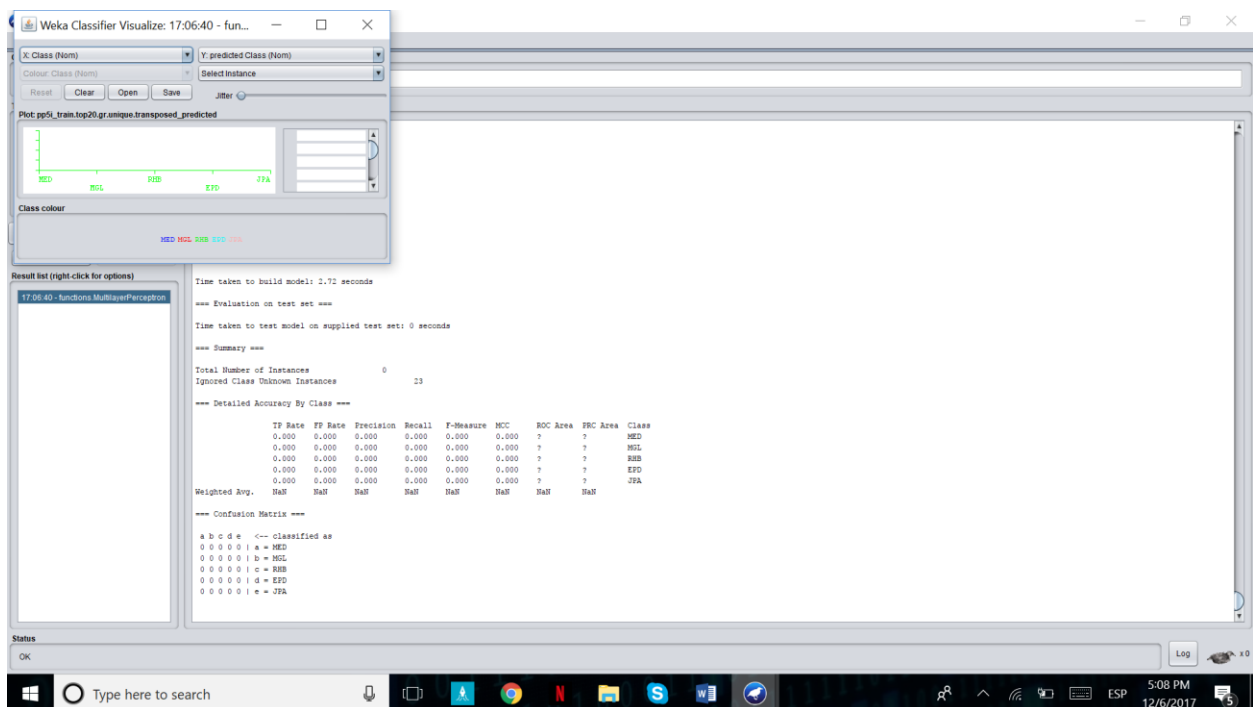


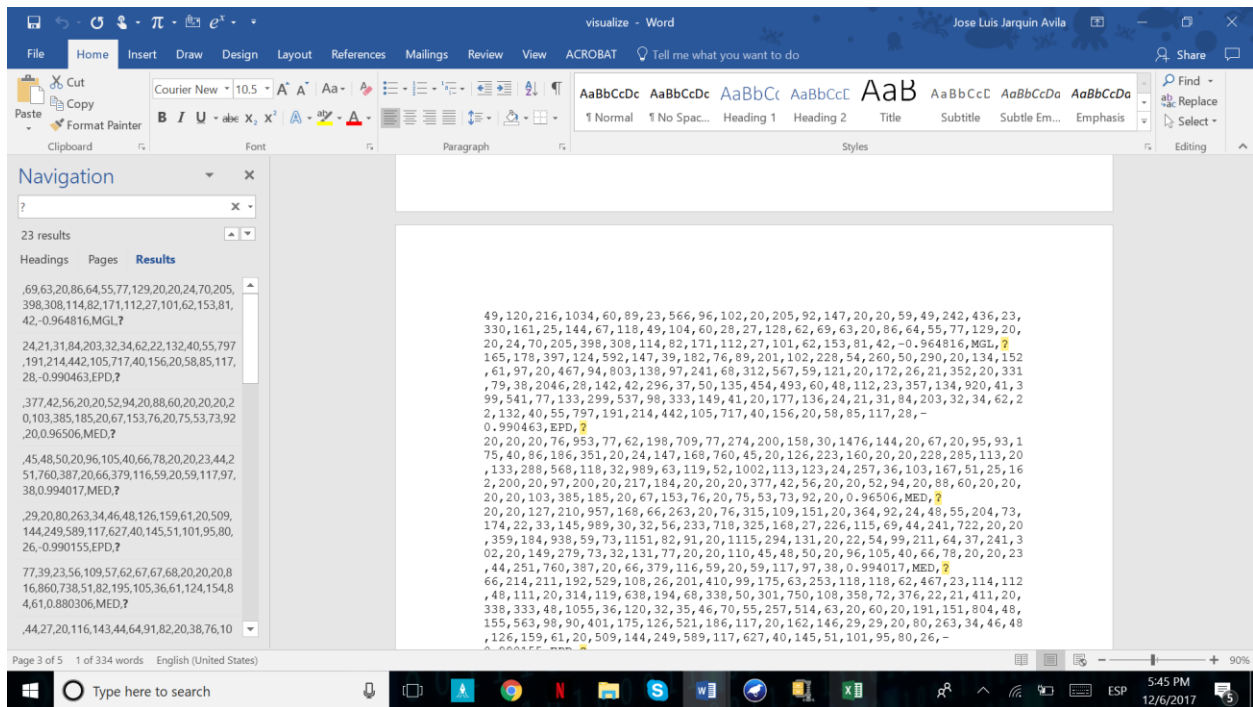*Figure 10 Weka's Visualize feature*
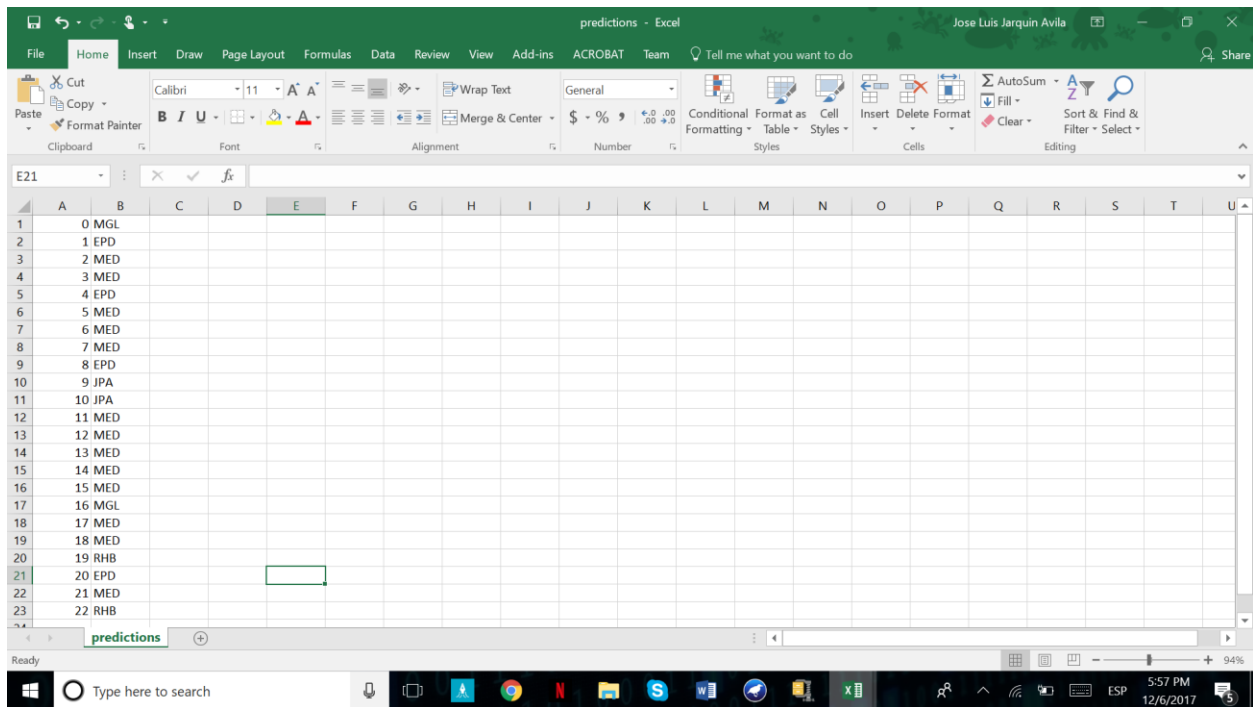
Figure 11 Adding the ? class

Figure 12 Getting the predicted classes in a column for easy transfer to the test file

**(5) (Recursively) Write a paper describing your paper.**

ZeroR is the simplest possible classifier. It selects the class that appears most often in the set.

IBk is Weka's implementation of the k-Nearest Neighbors classifier. 1-Nearest Neighbors (IB1) is also a simple classifier, by selecting the class of the closest training sample when given a test sample. 2-NN (IB2) and k-NN (IBk) generalizes this classifier by selecting the majority class of the N neighbors that lie closest to the test instance. In general, odd numbers of N are preferred to avoid ties.

Naïve Bayes is a simple classifier that operates on Bayes Theorem using the unrealistic (or at least generous) assumption that the classes in the data set are independent (hence Naïve). Bayes Theorem states that

$$p(C_k \mid \mathbf{x}) = \frac{p(C_k)\, p(\mathbf{x} \mid C_k)}{p(\mathbf{x})}$$

,

(courtesy Wikipedia) The term on the left-hand side is the posteriori probability, X is the evidence, and Ck are the classes. The Naïve Bayes classifier is thus constructed by making a frequency table of classes for each class, and taking the product of them divided by the likelihood of that class overall. A Laplace modifier is also usually added to ensure there is no class with zero probability. Test samples are classified by calculating the probability of the sample being a given class and the probability of the sample not being that class.

J48 is a decision-tree classifier. It works by recursively calculating the information gain of splitting on each class and constructing a tree. Once constructed, J48 classifies test samples by recursively walking down the tree and choosing to move to each child node that has the highest information gain (difference of entropy). This means that J48 uses the most predictive classes at every level of the tree to classify the instance.

Finally, we chose the top-20 unique genes from the classes, along with the Multilayer Perceptron, for our classifier. According to Weka, the multilayer perceptron is a classifier "that uses backpropagation to classify instances." Multilayer perceptrons have 3 layers, an input layer, an output layer, and a hidden layer. The train set is used to build the hidden layer, and then test inputs are fed into the input layer to classify the instances.

Below is a graph showing the incorrectly classified samples by class and gene subset size.
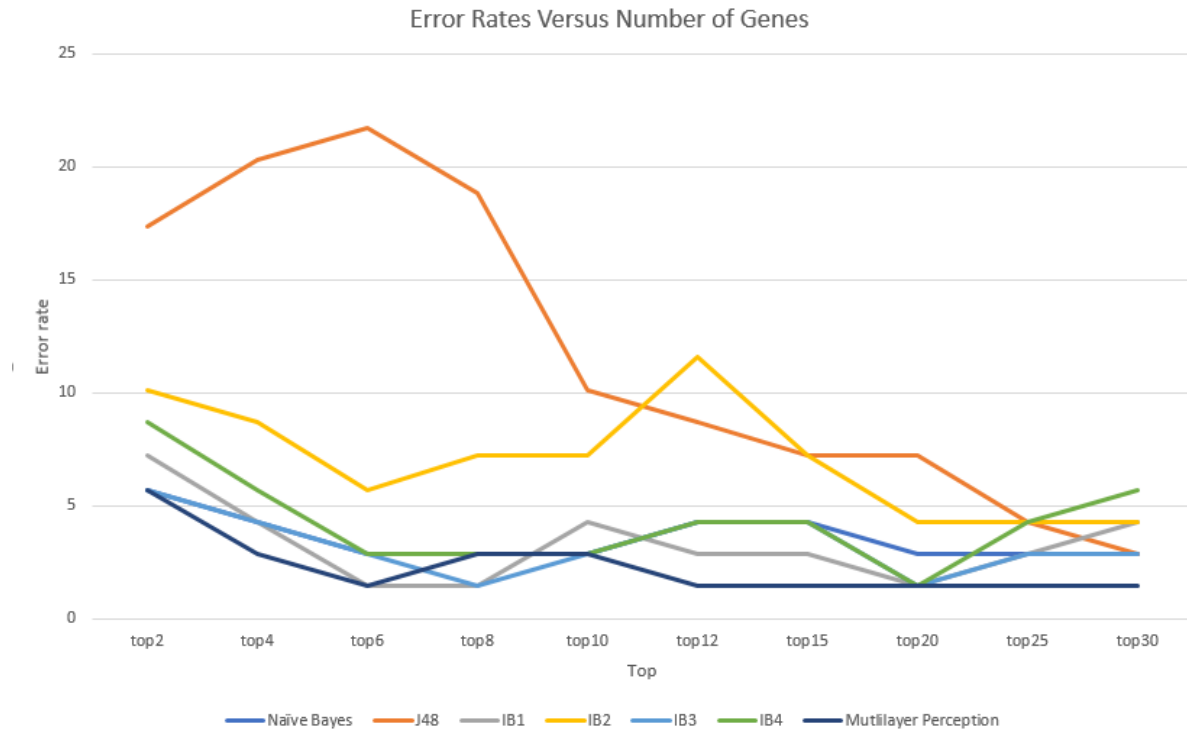
*Figure 13 Graph of each classifier's error rate against the top-n subsets generated.*