COMP 401 – Project in Biology and Computer Science

# Final Report

Supervised by: Dr. Thomas Bureau & Mr. Emilio Vello

Lea Collin, 260618407

Terms to know: unswitched case: case where training dataset has 185 instances and testing dataset has 908 instances (original way datasets were formed); switched case: case where training dataset has 908 instances and testing dataset has 185 instances; ranked case: case where training dataset has 908 instances with attributes ranked and testing dataset has 185 instances

## 1. Introduction

Climate change is becoming an ever-growing problem. The increase in global temperatures will lead to more drought events and some areas of the world are expected to experience reductions in their crop yields. For these areas, plant breeding programs and agricultural technology must be developed (Vello et al., 2015).

The use of different phenomics technologies in plants is a key element to improve our knowledge of the genotype-phenotype association of desired agricultural traits such as the response to water deficit (Vello et al., 2015). In the McGill Plant Phenomics Platform (MP3), research is being conducted in the hopes of better understanding plant response to water deficit, among others. The MP3 is instrumental in conducting this research. The MP3 has two main components, the 3D system and the Hts systems. The 3D system has the following cameras: 2 visible light, 2 neo-infrared, and 2 infrared, consisting of top and side views for each type of camera. This allows the phenomics lab to capture most of the characteristics of the plant. The system has a lifter with two different positions, high and low, and it can take images at single degree differences, ie 0 degrees, 1 degree, 2 degrees, etc. A conveyor belt transports the plants and the cabinets in which the plants are held have different positions that can rotate. There is a weighing system which allows for the tracking of automatic watering of the plant. The system can be programmed to keep the plant at a certain watering level or absolute weight value, eg. the plant will always receive 100 mL of water or the plant will receive water until it weighs 300 g (including the weight of the carrier). MP3 is located in the McGill greenhouse in the Stewart Biology building.

In this research, soybeans were studied because they are economically important. As the meat-packing industry is known to have very negative effects on the environment, soybeans have become even more important for the environmental sustainability of the future. Because soybean products are ingredients in many meat and dairy substitutes and soybeans are an exceptional source of essential nutrients, the shift from meat to soybeans may be essential (Arnarson et al., 2015). Learning more about soybeans would better allow us to effectively use it as a meat and dairy substitute and alleviate some negative environmental effects.

Machine learning and image analysis are important in this field so that we can better predict the effect that climate change will have on crop production and better prepare for it (Vello et al., 2015, Stavness et al., 2017). With the ability to instantly analyze different plants grown under different watering conditions, we will be able to get a better sense of what needs to be done in order to still meet the global demand for plant agriculture.

## 2. Background
### 2.1. Image Processing
All of the necessary image analysis was performed by Emilio Vello and the processed image data was stored in the MP3 soybean database. Figure 1 shows an example of the image analysis pipeline that was performed prior to this research.
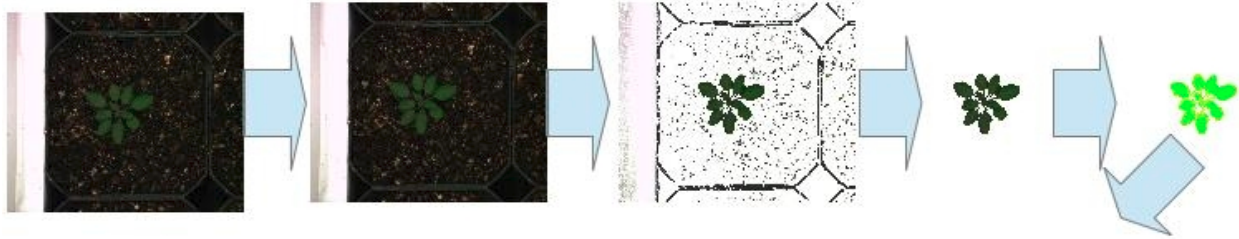
**Figure 1 –** Image analysis pipeline performed prior to research and stored in arabidopsis database; From "Image Analysis Pipeline", by Emilio Vello, 2017, http://mp3.biol.mcgill.ca/mcgill_mp3_analysis.html. Copyright 2017 by Emilio Vello. Reprinted with permission.

The final image in Figure 1 shows the image from which the data is taken and stored in the phenomics lab database.

### 2.2. Weka

Machine learning has an ever-growing important role in many fields, including image-based plant phenotyping. Weka is a collection of state-of-the-art machine learning algorithms (Frank and Witten, 1999). It was developed at the University of Waikato in New Zealand and the name stands for *Waikato Environment for Knowledge Analysis (*Frank and Witten, 1999). The system is written in Java and runs on almost any platform. Weka provides implementations of learning algorithms that one can easily apply to a dataset. All of the algorithms take their input in ARFF format which can be read from a file or generated by a database query (Frank and Witten, 1999). For this research, the database queries were saved as ARFF files and then read in by the program. One way of using Weka is to use learned models to generate predictions on new instances, which is the way that Weka was used for this research (Frank and Witten, 1999). Besides just training and building models based on a given training dataset, Weka also has several methods for "Attribute Selection". These different "Attribute Selection" methods read the training dataset and determine which attributes in the dataset are most important in predicting the developmental stage of the plant. One of these methods ranks the attributes by level of importance and these rankings can then be used while training and building the classification models. Weka is available through an easy-to-use graphical user interface or directly in code through the Java Weka library.

### 3. Methods
### 3.1. Software

During this research, a 'Weka Application' Java project was developed consisting of one class and several methods, that will be used by Dr. Thomas Bureau and Mr. Emilio Vello in their phenomics lab. For simplicity of building an initial program that Mr. Vello can later build upon, the developmental stage of the plant was based on its DAS (days after seeding). The different das were split into 5 categories, each containing 10 values and the developmental stages were then assigned to one of these 5 categories (Stage 1, Stage 2, Stage 3, Stage 4, or Stage 5). For example, a plant who was 34 days after seeding was considered to be in Stage 4.

This class provides an easy way to query the image data from the phenomics lab's soybean database, write the data to an ARFF file that is understandable to Weka methods, build and train classification models provided by the Java Weka library, and test these models on a testing dataset, which also comes from the soybean database. The jdbc drivers were used to be

able to write SQL queries within the class and the Java Weka library was used to directly import Weka methods. The Weka library provided all of the machine learning classification algorithms that were trained and tested within the class. The program simply reads the training data, builds a classification model based on this data and the name of the model inputted by the user (or builds the program's default model options). The program then runs on the testing dataset and calculates a precision of the model. All of the predictions of the different classification algorithms tested were outputted to a CSV file which also contained the barcode, das, and true stage of the plant. Another CSV is produced that stores the name of each algorithm tested and their respective precision value. The precision value was simply calculated as the number of instances that the model classifies correctly divided by the total number of instances within the dataset. These precision values matched those reported by Weka when performing the training and testing of the models directly within the Weka GUI. The CSV containing the stage prediction of each plant was then analyzed in R.

Both the Java code and R analysis code can be viewed at the end of this report.

### 3.2. Datasets and Tests

The training and testing sets were formed by querying the phenomics lab's soybean database. The public schema was used as were the dasplusev, imageev, imgobjectev, and soyidentification tables. These tables allowed the joining of the identification of the plant along with all of the important image characteristics that would be used to determine the developmental stage of the plant. The original training data only looked at the vis-side-1-0 camera view , line 1, and set 3. The original testing data also only looked at the vis-side-1-0 camera, lines 1, 2, and 3 from set 2 and lines 1 and 2 from set 3. It is important to only choose one camera for training and testing otherwise this will lead to non-comparable image views. The area of the plant from an aerial view will be calculated differently than the area of the plant from a side view. Therefore, it is important to be consistent with the type of camera when querying the data. The data selected additionally only included "well-watered" plants as opposed to also including "drought" plants. This is because the developmental stage of plants in these different conditions vary greatly. The term 'original' training data refers to the fact that this data is what was initially used to build the program (this will later on in the paper be referred to as the 'unswitched case', as described in 'terms to know' above).  The testing data was later on used as the training data and vice versa to see if there would be any improvement in precision (this will later on in the paper be referred to as the 'switched case', as described in 'terms to know' above).

### 3.3. Approach

The user was left with the option of inputting which classification algorithms they want to try. This was possible due to the "AbstractClassifier" class within the Weka library. The "AbstractClassifier" class would instantiate any classifier object simply based on the name of the classifier. This prevented the necessity of importing every classifier available from Weka regardless of user input. An array list of strings containing all of the possible classification algorithms was created in order to ensure that the user input was valid. This array list is also the 'default' set of classification algorithms that the program will run if the user decides not to input any algorithms on their own. Because the datasets are fairly small, it only took a few seconds for all of the models to be built and tested. The user was also allowed to decide whether they wanted to rank the attributes in the datasets to test if there would be any improvement in precision from doing so.
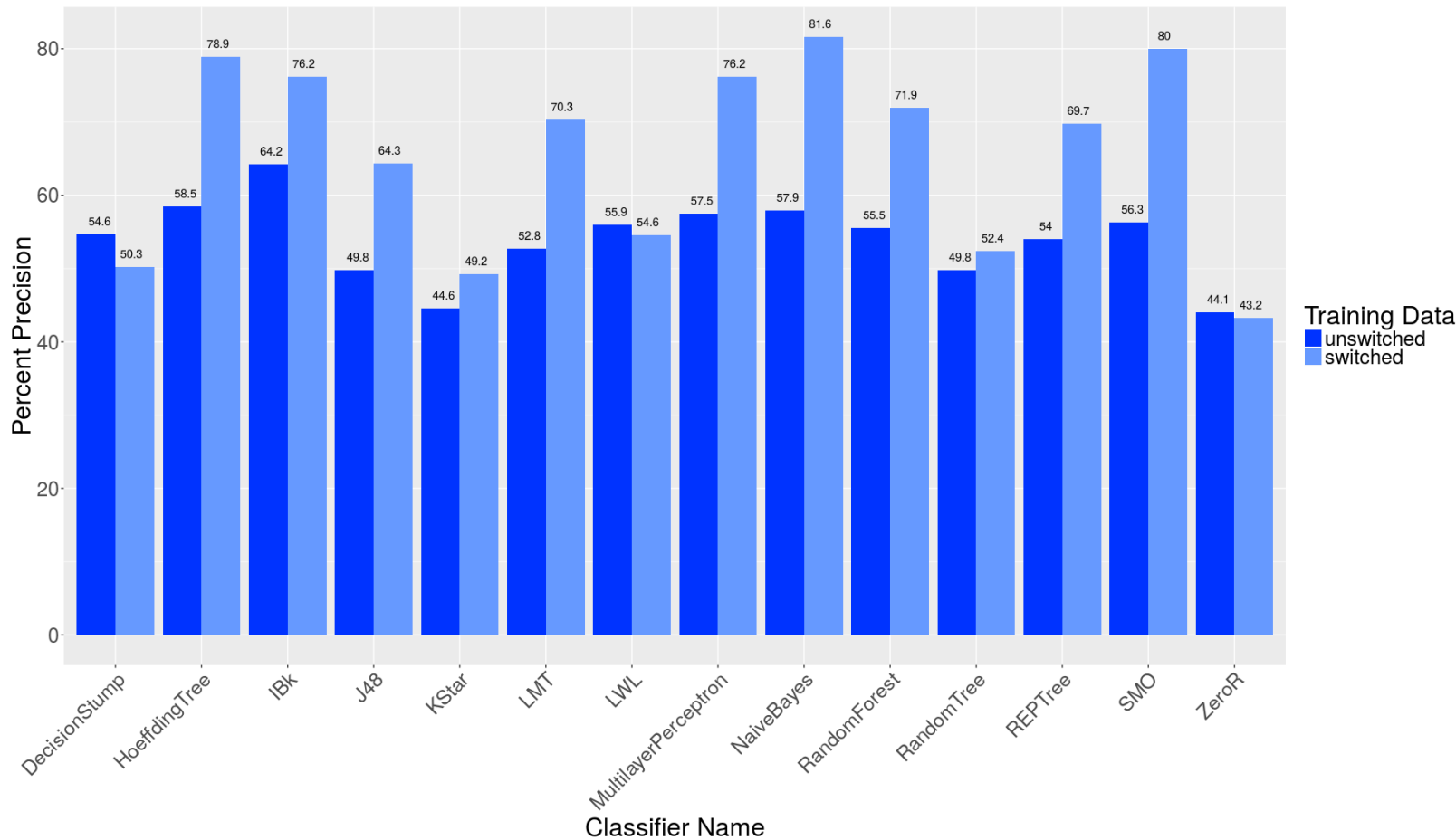
## 4. Results



**Figure 2** - All default algorithms tested with their precisions graphed side by side with the unswitched (dark blue) and switched (light blue) case

The precision values of all of the default algorithms are visible in Figure 2. The 'unswitched' bars in the figure refer to the 'original' organization of the training and testing data as described in 'terms to know' and Methods above. The 'switched' bars therefore refer to switching which dataset was used for training versus which dataset was used for testing. The following algorithms had better precision with the 'unswitched' datasets: ZeroR, LWL, and Decision Stump. The following algorithms had better precision with the 'switched' datasets: NaiveBayes, MultilayerPerceptron, SMO, IBk, KStar, HoeffdingTree, J48, LMT, RandomForest, RandomTree, and REPTree. Therefore, 11 algorithms performed better when the larger dataset was used for training and only 3 algorithms performed better when the smaller dataset was used for training.Overall, the best classification algorithm was NaiveBayes, in the 'switched' case, with 81.6% precision.
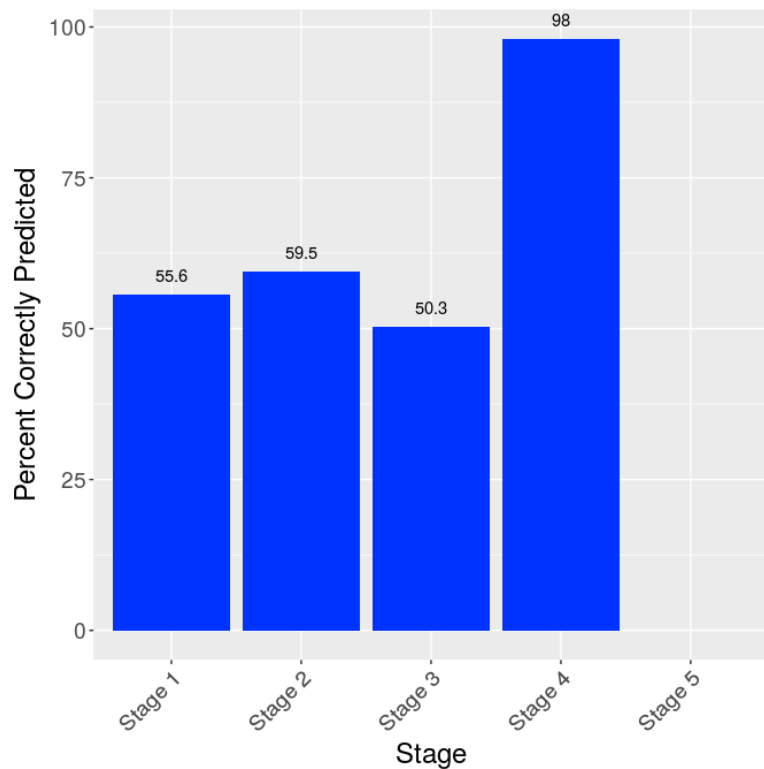
**Figure 3A –** Number of times each stage was correctly predicted as a percentage of the total number of times each stage was predicted; Sum of number of times each classifier correctly predicted the stage was calculated and divided by total number of times each stage was predicted at all; Average over all default algorithms except for ZeroR
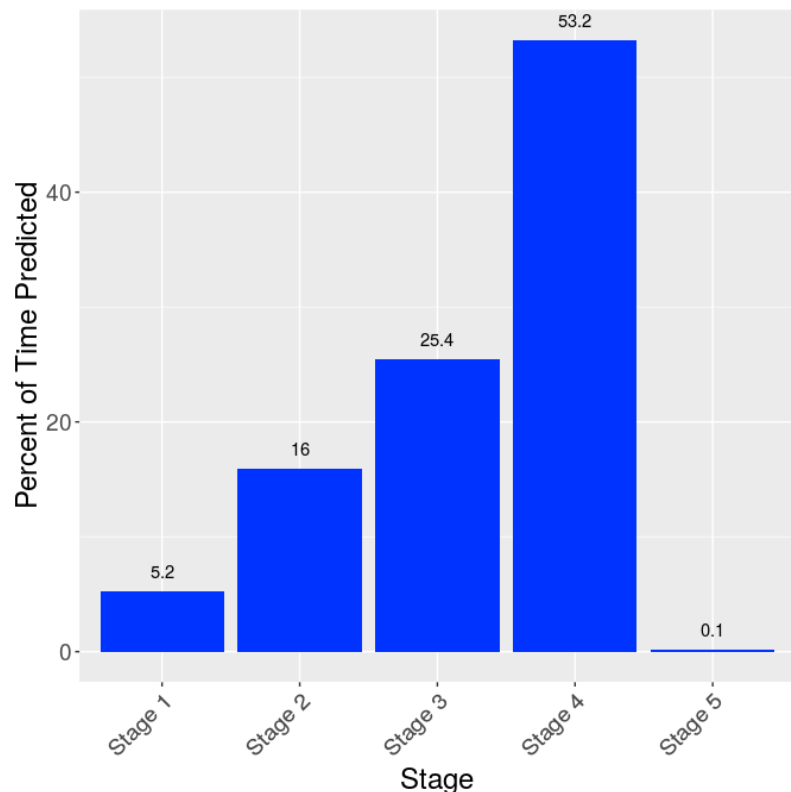
**Figure 3B –** Frequency that each stage was predicted; Number of times each stage was predicted by each classifier divided by total number of prediction made; Average over all default algorithms except for ZeroR

If we look at just the 'switched' case, ie with the larger dataset used for training and the smaller used for testing, Figure 3A shows the number of times that each stage was correctly predicted, as a percentage of the total number of times each stage was predicted. The sum of the number of times each classifier correctly predicted the stage was calculated and then divided by the total number of times each stage was predicted at all. Note that in these calculations, the ZeroR algorithm was ignored as this algorithm simply predicts the most common stage in the training dataset. For example, in the original training set, Stage 4 was the most common stage so the ZeroR algorithm predicted Stage 4 for each instance in the testing set, therefore including this algorithm would skew the data too much. The frequency that each stage was predicted for this case was also calculated. This frequency was calculated by counting the number of times each stage was predicted by all of the algorithms and then dividing by the total number of predictions made by all of the algorithms, except ZeroR. The frequencies resulted in the following breakdown: Stage 1 – 5.2%, Stage 2 – 16%, Stage 3 – 25.4%, Stage 4 – 53.2%, Stage 5 – 0.1%,which is shown in Figure 3B.

When looking at the 'unswitched' case, the stage precision breakdown we saw was the following: Stage 1 – 63.5%, Stage 2 – 73%, Stage 3 – 29%, Stage 4 – 79.5%, Stage 5 – 0%. The frequency that each stage was predicted for this case had the following breakdown: Stage 1 – 7.2%, Stage 2 – 32.2%, Stage 3 – 17.1%, Stage 4 – 43.5%, Stage 5 – 0%.
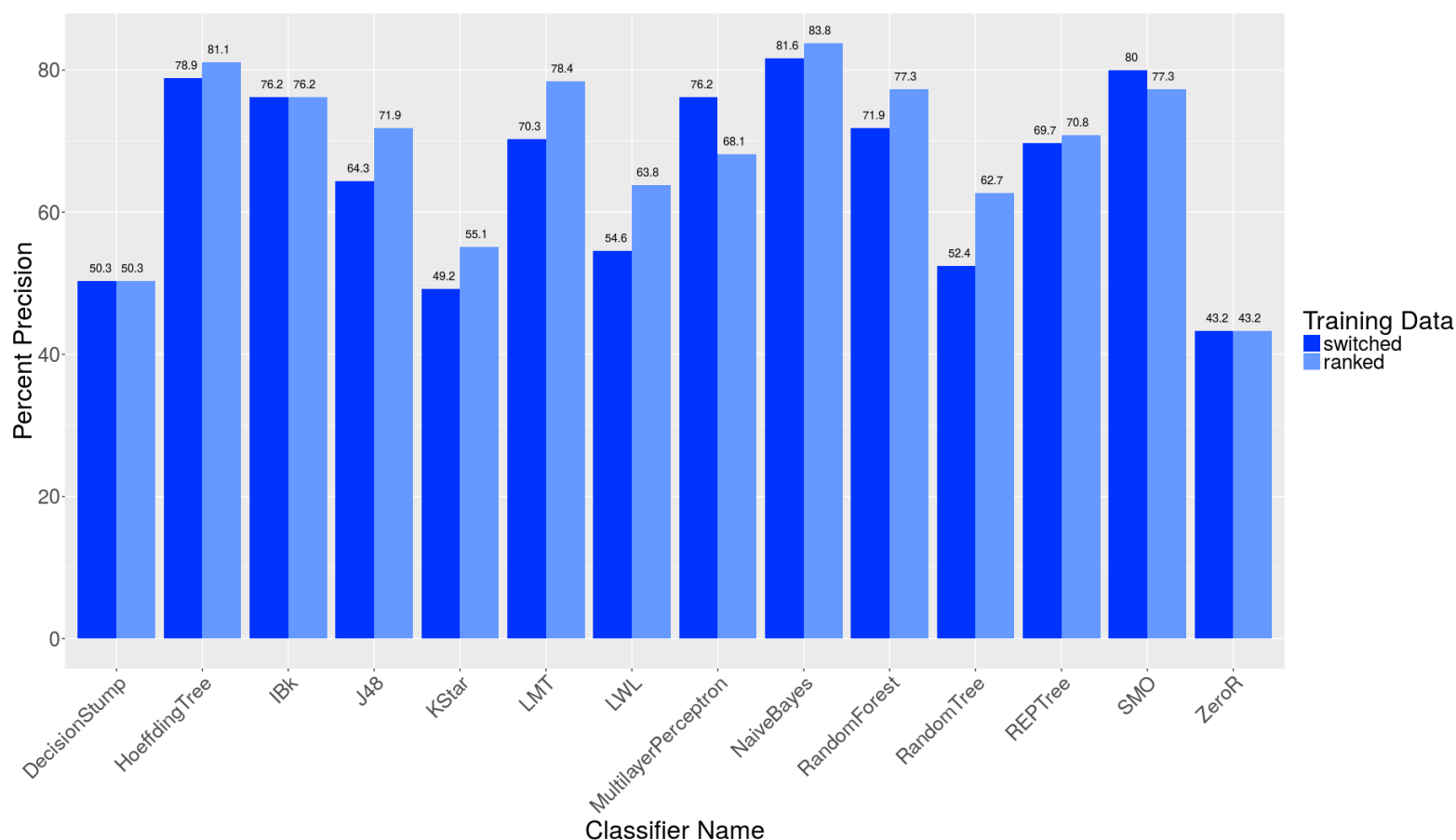
**Figure 4** - All defaut algorithms tested with their precisions graphed side by side with the ranked (dark blue) and switched (light blue) case

Figure 4 shows the different precision values for the 'switched' case and the 'ranked' case. The 'ranked' case took the 'switched' data and first ranked the attributes before training and building the classification models. The same default set of models were built and tested and the precision values were compared to those achieved when the ranking algorithm was not run. The algorithms that performed better when the attributes were ranked first were: NaiveBayes, KStar, LWL, HoeffdingTree, J48, LMT, RandomForest, RandomTree, and REPTree. The algorithms that performed better when the attributes were not ranked first were: MultilayerPerceptron and SMO. The algorithms that showed no difference in precision when the attributes were ranked or unranked were: ZeroR, IBk, DecisionStump. Therefore, 9 algorithms performed better when the attributes were ranked, 2 performed worse, and 3 performed the same. The best classification algorithm after ranking the attributes was the NaiveBayes algorithm with a precision of 83.8%, a 2.2% increase in precision from NaiveBayes in just the 'switched' case with no attribute ranking.

When looking at the 'ranked' case, the stage precision breakdown we saw was the following: Stage 1 – 79.6%, Stage 2 – 81.5%, Stage 3 – 42.1%, Stage 4 – 99.5%, Stage 5 – 0%. The frequency that each stage was predicted for this case had the following breakdown: Stage 1 – 6.6%, Stage 2 – 19.8%, Stage 3 – 18.3%, Stage 4 – 55.3%, Stage 5 – 0%.

**5. Discussion**

From Figure 3A (the 'switched' case), we see that Stage 4 is predicted drastically more precisely than any other stage. From Figure 3B, we also see that Stage 4 was predicted the most often. Looking at a breakdown of the number of instances belonging to each stage in the dataset we have: (Stage 1: 32, Stage 2: 163, Stage 3: 280, Stage 4: 400, Stage 5: 33). The combination of Stage 4 being the most prevalent stage in the set used for training and the fact that it was predicted the most often, gives a possible reason that this stage was predicted most accurately. Additionally, Stage 4 plants are larger since they are more days after seeding. These plants may be more developed and more defined than smaller plants which have less information to be gained from them. This gives another possible reason that Stage 4 plants were predicted most precisely. Stage 5 was never predicted correctly because in this case the testing set had zero plants in Stage 5.

In the 'unswitched' case, we see again that Stage 4 was correctly predicted the most whereas Stage 3 was correctly predicted the least often (excluding Stage 5). We also see that Stage 4 was predicted the highest number of times but we also see that Stage 3 was not predicted the fewest number of times, Stage 1 was (again excluding Stage 5). Looking at a breakdown of the number of instances belonging to each stage in the dataset we have: (Stage 1: 9, Stage 2: 36, Stage 3: 60, Stage 4: 80, Stage 5: 0). We can see from this that there were no instances in the dataset that belonged to Stage 5 and this is why Stage 5 was never predicted nor correctly predicted, the classification algorithms had nothing to learn from. Similar to the 'switched' case, Stage 4 was most likely predicted most accurately due to the fact that there is more information to be gained from plants in this stage than plants in earlier stages. However, there is not yet a good explanation why Stage 3 was correctly predicted the least often instead of Stage 1.

From Figure 4 (the 'ranked' case) we see that ranking the attributes before building and training the models improved (or at least did not affect) the precision of most algorithms. For some algorithms, it is obvious why there were no improvements in precision when the attributes were ranked. The ZeroR algorithm, for example, always predicts the value that is most common in the training dataset. Therefore, it does not matter how well the data is prepared before building the model, the ZeroR algorithm will always predict the same value. Most of the algorithms experienced an improvement in precision as a result of ranking the attributes. This makes sense as the ranking tells the algorithms which attributes are most important to analyze to determine the developmental stage of the plant and therefore these attributes receive more weight. It seems reasonable that in the worst-case, this ranking would just be extra computational work that added no benefit to precision and the precision would simply be the same as if we did not rank first, which is what we saw for three of the algorithms tested. We once again saw that Stage 4 was the best predicted stage when the attributes were ranked and we can again guess that this is due to the plants in Stage 4 being larger and better defined than the plants in earlier stages. What is interesting, however, is that we see a 24% increase in precision in predicting Stage 1, a 22% increase in precision in predicting Stage 2, and an 8.2% decrease in predicting Stage 3. As just stated, the increase in precision for certain stages is most likely due to the fact that the more important attributes have more weight being assigned to them. However, there is not a good explanation yet for why Stage 3 would have experienced a decrease in precision compared to when the attributes were unranked.

## 6. Conclusion

In this research, a program was developed as a way of automating the querying of any database from the McGill phenomics lab, reading this data, using the training data to build machine learning classification models that would predict the developmental stage of a plant based on its image data, and testing these models on a new test set of plant image data. Based on the program, the NaiveBayes algorithm is most precise in correctly predicting the developmental stage based on a new set of data. This research also showed that a larger training set leads to better models, which seems reasonable and makes sense. As we saw here, the initial training set had 185 instances while the testing set had 908 instances. When the smaller dataset is used for training, the highest precision we achieved from any algorithm was only 64.2% (from the IBk algorithm). However, when the larger dataset was used for training and the smaller dataset was used for testing, ie the 'switched' case, we were able to achieve 81.6% precision (from the NaiveBayes algorithm). We were then able to reach an even slightly higher precision by looking at the 'switched' case and ranking the attributes before building the models. With this method, we were able to achieve 83.8% precision (again from the NaiveBayes algorithm).

We saw from the R analysis that Stage 5 was predicted the least precisely both in the 'switched' and 'unswitched' cases and we safely concluded that this is due to the small number of plants that were in Stage 5 in both datasets, relative to the other stages. We saw that Stage 4 was predicted the most precisely in both the 'switched' and 'unswitched' cases. We concluded that this was due to the fact that plants in this stage were larger and more developed, thereby being more defined than smaller plants which would have less information available to learn about them.

This program provides an easy way to test many different algorithms quickly and efficiently connected to the MP3 database. A user can easily build and train models with any algorithm that exists within the Weka library. Additionally, the program also provides direct querying from the MP3 database, though it can be easily modified to query another database. A possible future opportunity is to see if different datasets lead to even higher precision. For example, in this research we only looked at the vis-side-1-0 camera but it is very possible to choose a different camera and angle to see if there would be an improvement in precision results. Though I implemented the option for the user to rank the attributes in the dataset by which attributes are most important in determining the developmental stage, I did not have the chance to implement the other attribute selection options that Weka provides. For example, Weka has different attribute selection options that will actually eliminate certain attributes and only keep the attributes that are most important in determining the developmental stage of the plant. A further step for this research would be to implement these different attribute selection options that the user can choose from to see if there would be any improvement in precision, similar to the improvement that was observed when the attribute ranking was performed.

References

Arnarson, Atli. (2015, January 17). *Soybeans 101: Nutritions Facts and Health Effects*. Retrieved from https://www.healthline.com/nutrition/foods/soybeans.

Frank, Eibe and Witten, Ian H. (1999). *Data Mining: Practical Machine Learning Tools and Techniques*. London: Morgan Kaufmann Publishers.

Stavness, Ian and Ubbens, Jordan R. (2017). "Deep Plant Phenomics: A Deep Learning Platform for Complex Plant Phenotyping Tasks". *Frontiers in Plant Science, 8,* 11.

Vello E., Diallo Oury, A., Tomita, A., and Bureau T.E., (2015). "A Comprehensive Approach to Assess Arabidopsis Survival Phenotype in Water-Limited Condition Using a Non-invasive High-Throughput Phenomics Platform". *Frontiers in Plant Science, 6,* 10.

Vello, Emilio. (2017). *Image Analysis Pipeline.* Retrieved from http://mp3.biol.mcgill.ca/mcgill_mp3_analysis.html

Vello, Emilio. (2017). *The Platform in a Short Presentation.* Retrieved from http://mp3.biol.mcgill.ca/mcgill_mp3_summary.html

# SUPPLEMENTARY MATERIALS

```java
/*
 * @author Lea Collin
 */
import weka.core.Instance;
import weka.core.Instances;
import weka.classifiers.Classifier;
import weka.classifiers.AbstractClassifier;
import weka.classifiers.meta.FilteredClassifier;
import weka.core.converters.ArffLoader;
import weka.filters.unsupervised.attribute.Remove;
import weka.attributeSelection.*;
import java.io.File;
import java.io.FileNotFoundException;
import java.io.PrintWriter;
import java.io.IOException;
import java.util.ArrayList;
import java.util.Scanner;
import org.postgresql.util.PSQLException;
import java.sql.*;

public class Driver {

    static Scanner sc = new Scanner(System.in);

    public static void main(String [] args) throws Exception {

        //private database information
        String dbConfig = args[0];

        //ArrayList will contain all the possible algorithms that the user
  can input
        ArrayList<String> possibleClassifiers = new ArrayList<String>();
        possibleClassifiers.add("ZeroR");
        possibleClassifiers.add("NaiveBayes");
        possibleClassifiers.add("MultilayerPerceptron");
        possibleClassifiers.add("SMO");
        possibleClassifiers.add("IBk");
        possibleClassifiers.add("KStar");
        possibleClassifiers.add("LWL");
        possibleClassifiers.add("DecisionStump");
        possibleClassifiers.add("HoeffdingTree");
        possibleClassifiers.add("J48");
        possibleClassifiers.add("LMT");
        possibleClassifiers.add("RandomForest");
        possibleClassifiers.add("RandomTree");
        possibleClassifiers.add("REPTree");

        System.out.println("Please enter the directory name of where you
```

```
     would like to store all program outputs:");
48
49            //add control to check for valid directory
50            String outputDir = setOutputDirectory();
51
52            System.out.println();
53            System.out.println("Please enter the name of the file you'd like to
   store the TRAINING data. Please end the file name in '.arff'");
54            String trainingFile = setFileName(".arff", outputDir);
55            //String trainingFile = outputDir + "trainingData.arff";
56            System.out.println();
57
58            System.out.println("Please enter the name of the file you'd like to
   store the TESTING data. Please end the file name in '.arff'");
59            String testingFile = setFileName(".arff", outputDir);
60            System.out.println();
61
62            //try to connect to database given username and password, user
   prompted to enter username and password again if connection is unsuccessful
63            boolean successfulConnection = false;
64            while(!successfulConnection) {
65                System.out.println("Please enter your database username:");
66                String dbUsr = sc.next();
67                System.out.println("Password:");
68                String dbPwd = sc.next();
69                sc.nextLine();
70            try {
71                connectToDatabase(dbUsr, dbPwd, dbConfig, testingFile,
   trainingFile);
72
73                successfulConnection = true;
74            }catch (PSQLException s){
75                System.out.println("Username or password was incorrect.
   Please try again.");
76                }
77            }
78
79            boolean rankAttributes = false;
80            boolean isValid = false;
81            while(!isValid) {
82                System.out.println("Would you like to rank the attributes? (Y/n),
   (Could lead to better results)");
83                String answer = sc.next();
84                if(answer.equals("Y") || answer.equals("y")) {
85                    isValid = true;
86                    rankAttributes = true;
87                }
88                else if(answer.equals("N") || answer.equals("n")){
```

```
89                    isValid = true;
90                    rankAttributes = false;
91                }
92                else {
93                    System.out.println("Could not understand input. Please enter
   a name again. \n");
94                }
95            }
96
97        //getting user input for classifier names, checking if input is valid
98        String [] classifiers = null;
99
100        boolean validClassifier = false;
101        while(!validClassifier){
102            System.out.println("Do you want to run the default set of
   classification algorithms? (Y/n)");
103            String answer = sc.next();
104            if(answer.equals("Y") || answer.equals("y")) {
105                classifiers = new String[possibleClassifiers.size()];
106                for(int i = 0; i < classifiers.length; i++) {
107                    classifiers[i] = possibleClassifiers.get(i);
108                }
109                validClassifier = true;
110            }
111            else if(answer.equals("N") || answer.equals("n")){
112                System.out.println("Please enter the names of the classifiers
   you'd like to test, separated by a single space.");
113                sc.nextLine();
114                String classifierInput = sc.nextLine();
115                classifiers = classifierInput.split("\\s+");
116
117                for(int i = 0; i < classifiers.length; i++) {
118                    if(!possibleClassifiers.contains(classifiers[i])) {
119                        System.out.println(classifiers[i] + " is not a valid
   classifier name.");
120                        System.out.println("Please try again.");
121                        System.out.println();
122                    }
123
124                    if(i == classifiers.length - 1 &&
   possibleClassifiers.contains(classifiers[i])) {
125                        validClassifier = true;
126                    }
127                }
128            }
129            else {
130                System.out.println("Could not understand input.");
131            }
```

```java
132
133          }
134
135          System.out.println();
136          System.out.println("Please enter the name of the file you would like
    to store all of the predictions. Please end the file in '.csv'");
137          String predictionFile = setFileName(".csv", outputDir);
138
139          System.out.println();
140          System.out.println("Finally, please enter the name of the file you
    would like to store the precision of each algorithm you are testing. "
141                    + "Please end the file in '.csv'");
142          String precisionFile = setFileName(".csv", outputDir);
143          System.out.println();
144
145          sc.close();
146
147          //reading the files and getting all the instances of each one
148          Instances instancesTrain = fileReader(trainingFile);
149          Instances instancesTest = fileReader(testingFile);
150
151          //what attribute do we want to predict
152          String classAttribute = "Stage";
153
154          //these are the attributes to not be included in the classifier, das
    is a giveaway of the stage, barcode is irrelevant
155          String [] attributesToRemove = {"barcode", "das"};
156          //returns the indices of the attributes to be ignored
157          String indicesToRemove = removeAttribute(instancesTrain,
    attributesToRemove);
158
159          double maxPrecision = 0.0;
160
161          String bestMethod = "";
162
163          //run all the different classifiers
164          Double [] precisions = predict(instancesTrain, instancesTest,
    classifiers, indicesToRemove, classAttribute, predictionFile,
    rankAttributes);
165
166          //writing precision values to a csv
167                  File output = new File(precisionFile);
168                  PrintWriter pw = new PrintWriter(output);
169                  StringBuilder sb = new StringBuilder();
170                  sb.append("Algorithm");
171                  sb.append(',');
172                  sb.append("Precision");
173                  sb.append("\n");
```

```java
174
175        //find best algorithm and write to file
176        for(int i = 0; i < classifiers.length; i++) {
177            if(precisions[i] > maxPrecision) {
178                maxPrecision = precisions[i];
179                bestMethod = classifiers[i];
180            }
181        }
182        //simply printing out for the user which algorithm was best and its
    precision
183        System.out.println("Best Algorithm: " + bestMethod + " with
    precision: " + maxPrecision);
184
185        pw.write(sb.toString());
186        pw.close();
187    }
188
189    public static boolean validArff(String file) {
190        //string must contain and end in .arff to be a valid arff file
191        return (file.contains(".arff") && file.indexOf(".arff") ==
    file.length() - 5);
192    }
193
194    public static boolean validCsv(String file) {
195        return (file.contains(".csv") && file.indexOf(".csv") == file.length
    () - 4);
196    }
197
198    public static boolean fileExists(String file) {
199        File newFile = new File(file);
200        return newFile.exists();
201    }
202
203    public static String setOutputDirectory() {
204        String dirName = "";
205        boolean isValid = false;
206        while(!isValid) {
207            dirName = sc.next();
208            File dir = new File(dirName);
209
210            if(dir.exists() && (dirName.charAt(dirName.length()-1) == '/')) {
211                isValid = true;
212            }
213            else {
214                System.out.println("Sorry that was an invalid directory.
    Please try again.");
215            }
216        }
```

```java
217
218            return dirName;
219        }
220
221        public static String setFileName(String fileType, String outputDir) {
222            String outputFile = "";
223            boolean isValid = false;
224            while(!isValid) {
225                outputFile = outputDir + sc.next();
226                if(fileType.equals(".arff")) {
227                    if(validArff(outputFile) == isValid) {
228                        System.out.println("Sorry, the file you entered does not
   end in '.arff'. Please try again.");
229                        continue;
230                    }
231                }
232                if(fileType.equals(".csv")) {
233                    if(validCsv(outputFile) == isValid) {
234                        System.out.println("Sorry, the file you entered does not
   end in '.csv'. Please try again.");
235                        continue;
236                    }
237                }
238                if(fileExists(outputFile)) {
239                    System.out.println("This file already exists in this
   directory. Do you want to overwrite it? (Y/n)?");
240                    String answer = sc.next();
241                    if(answer.equals("Y") || answer.equals("y")) {
242                        isValid = true;
243                    }
244                    else if(answer.equals("N") || answer.equals("n")){
245                        System.out.println("Please enter another name.");
246                    }
247                    else {
248                        System.out.println("Could not understand input. Please
   enter a name again.");
249                    }
250                }
251                else {
252                    isValid = true;
253                }
254            }
255            return outputFile;
256        }
257
258        public static Instances fileReader(String input) throws IOException {
259
260            File inputFile = new File(input);
```

```java
261        ArffLoader atf = new ArffLoader();
262        atf.setFile(inputFile);
263        Instances data = atf.getDataSet();
264
265        return data;
266
267    }
268
269    public static String removeAttribute(Instances data, String []
   attributes) throws Exception{
270
271        String [] options = new String[2];
272        options[0] = "-R";
273
274        String indices = "";
275
276        //getting all the indices of attributes we want to later remove (or
   ignore)
277        for(int i = 0; i < attributes.length; i++) {
278            int index = (data.attribute(attributes[i])).index() + 1;
279            indices += index;
280            if(i != attributes.length-1) {
281                indices += ",";
282            }
283        }
284
285        return indices;
286    }
287
288    public static Double [] predict(Instances train, Instances test, String
   [] classifierNames,
289        String indicesToRemove, String classAttribute, String outputFile,
   boolean rankAttributes) throws Exception {
290
291        //removing attributes we don't want to include such as das and
   barcode
292                Remove rm = new Remove();
293                rm.setAttributeIndices(indicesToRemove);
294
295        //running an attribute selector if user chose to rank attributes
296                if(rankAttributes) {
297                    Instances [] reduced = attributeSelector(train, test);
298
299                    train = reduced[0];
300
301                    test = reduced[1];
302                }
303
```

```
304         //set the Class (what we want to predict)
305         test.setClass(test.attribute(classAttribute));
306
307         //setting the train class index to be the same as the testing class
    index
308         train.setClassIndex(test.classIndex());
309
310         //going to make an array of classifiers
311         Classifier [] classifiers = new Classifier [classifierNames.length];
312
313         //need somewhere to store the precision of each classifier +
    initializing the array
314         Double precision [] = new Double [classifiers.length];
315         for(int i = 0; i < precision.length; i++) {
316             precision[i] = 0.0;
317         }
318
319         //creating filtered versions of each classifier (removing das and
    barcode)
320         for(int i = 0; i < classifiers.length; i++) {
321             Classifier temp = AbstractClassifier.forName(classifierNames[i],
    null);
322
323             FilteredClassifier fc = new FilteredClassifier();
324             fc.setFilter(rm);
325             fc.setClassifier(temp);
326
327             classifiers[i] = fc;
328         }
329
330         //building the models for each classifier
331         for(int i = 0; i < classifiers.length; i++) {
332             classifiers[i].buildClassifier(train);
333         }
334
335         //writing the header to the output csv
336         File output = new File(outputFile);
337         PrintWriter pw = new PrintWriter(output);
338         StringBuilder sb = new StringBuilder();
339         sb.append("Barcode");
340         sb.append(',');
341         sb.append("Das");
342         sb.append(",");
343         sb.append("Actual");
344         sb.append(',');
345
346         //making the names of the classifiers part of the header, will
    indicate the stage that classifier has predicted for a particular plant
```

```
347            for(int i = 0; i < classifierNames.length; i++) {
348                sb.append(classifierNames[i]);
349                if(i != classifierNames.length - 1) {
350                    sb.append(",");
351                }
352            }
353            sb.append("\n");
354
355            double numInst = test.numInstances();
356
357            //actually running the classifier
358            for(int i = 0; i < numInst; i++){
359
360                Instance current = test.instance(i);
361
362                //will set the Stage of this 'temp' instance to be the predicted
   value to then compare it to the actual value of 'current'
363                Instance temp = (Instance)current.copy();
364
365                //attributes are given as array positions, getting the string
   value
366                String actualVal = current.stringValue(test.classIndex());
367
368                sb.append((int) current.value(test.attribute("barcode")));
369                sb.append(',');
370                sb.append((int) current.value(test.attribute("das")));
371                sb.append(",");
372                sb.append(actualVal);
373                sb.append(',');
374
375              for(int j = 0; j < classifiers.length; j++) {
376                //getting the predicted value of the class attribute of this
   instance
377                double predicted = classifiers[j].classifyInstance(test.instance
   (i));
378
379                //setting this value to the temp class attribute
380                temp.setValue(test.classIndex(), predicted);
381
382                //getting the string value
383                String predictedVal = temp.stringValue(temp.classIndex());
384
385                //comparing predicted with actual value
386                if(predictedVal.equals(actualVal)) {
387                    precision[j]++;
388                }
389
390                sb.append(predictedVal);
```

```
391
392                if(j != classifiers.length - 1) {
393                    sb.append(",");
394                }
395
396                else if( j == classifiers.length - 1) {
397                    sb.append('\n');
398                }
399
400            }
401        }
402
403        for(int i = 0; i < precision.length; i++) {
404            precision[i] = 100*precision[i]/numInst;
405        }
406
407        sb.append('\n');
408        pw.write(sb.toString());
409        pw.close();
410
411        return precision;
412    }
413
414    private static void connectToDatabase(String usrDB, String passwordDB,
    String conDB, String trainName, String testName) throws SQLException,
    FileNotFoundException {
415
416        File trainingOutput = new File(trainName);
417        PrintWriter trainingPw = new PrintWriter(trainingOutput);
418        StringBuilder sb = new StringBuilder();
419
420        File testingOutput = new File(testName);
421        PrintWriter testingPw = new PrintWriter(testingOutput);
422
423        //creating the header for the arff file
424        sb.append("@relation databasetraining" + "\n" +  "\n" + "@attribute
    barcode numeric"  + "\n" + "@attribute area numeric" + "\n" + "@attribute
    perimeter numeric" + "\n" +
425            "@attribute circularity numeric" + "\n" + "@attribute compactness
    numeric" + "\n" + "@attribute major numeric" + "\n" +
426            "@attribute minor numeric" + "\n" + "@attribute eccentricity numeric"
    + "\n" + "@attribute hisgreypeak numeric" + "\n" +
427            "@attribute q1grey numeric" + "\n" + "@attribute q2grey numeric" +
    "\n" + "@attribute q3grey numeric" + "\n" +
428            "@attribute q1r numeric" + "\n" + "@attribute q2r numeric" + "\n" +
    "@attribute q3r numeric" + "\n" + "@attribute q1g numeric" + "\n" +
429            "@attribute q2g numeric" + "\n" + "@attribute q3g numeric" + "\n" +
    "@attribute q1b numeric" + "\n" + "@attribute q2b numeric" + "\n" +
```

```
430          "@attribute q3b numeric" + "\n" + "@attribute das numeric" + "\n"
    +"@attribute Stage {'Stage 1','Stage 2','Stage 3','Stage 4', 'Stage 5'}" +
    "\n" + "\n" + "@data" + "\n");
431
432          trainingPw.write(sb.toString());
433          testingPw.write(sb.toString());
434
435          try {
436              Class.forName("org.postgresql.Driver");
437              Connection conn = DriverManager.getConnection(conDB, usrDB,
    passwordDB);
438
439              String trainingSql = "SELECT s.barcode, o.area, "
440                      + "o.perimeter, o.circularity, o.compactness, "
441                      + "o.major, o.minor, o.eccentricity, o.hisgreypeak, "
442                      + "o.q1grey, o.q2grey, o.q3grey, "
443                      + "o.q1r, o.q2r, o.q3r, "
444                      + "o.q1g, o.q2g, o.q3g, "
445                      + "o.q1b, o.q2b, o.q3b, d.das, "
446                      + "CASE WHEN ( d.das <= 10 ) THEN 'Stage 1' "
447                      + "WHEN ( d.das > 10 AND d.das <= 20 ) THEN 'Stage 2' "
448                      + "WHEN ( d.das > 20 AND d.das <= 30 ) THEN 'Stage 3' "
449                      + "WHEN ( d.das > 30 AND d.das <= 40 ) THEN 'Stage 4' "
450                      + "WHEN ( d.das > 40 AND d.das <= 50 ) THEN 'Stage 5' "
451                      + "ELSE 'Stage 6' END Stage "
452                      + "FROM imageev AS i, imgobjectev AS o, soyidentification
    AS s, dasplusev AS d "
453                      + "WHERE i.assayid = o.assayid "
454                      + "AND i.imgid = o.imgid "
455                      + "AND s.barcode = ( CAST( i.barcode AS INTEGER ) ) "
456                      + "AND i.assayid = d.assayid AND i.fdate = d.fdate "
457                      + "AND i.set = d.set "
458                      + "AND s.line = 1 "
459                      + "AND i.camera = 'vis-side-1-0' "
460                      + "AND i.set = '3'";
461
462              String testingSql =  "SELECT s.barcode, o.area, "
463                      + "o.perimeter, o.circularity, o.compactness, "
464                      + "o.major, o.minor, o.eccentricity, o.hisgreypeak, "
465                      + "o.q1grey, o.q2grey, o.q3grey, "
466                      + "o.q1r, o.q2r, o.q3r, "
467                      + "o.q1g, o.q2g, o.q3g, "
468                      + "o.q1b, o.q2b, o.q3b, d.das, "
469                      + "CASE WHEN ( d.das <= 10 ) THEN 'Stage 1' "
470                      + "WHEN ( d.das > 10 AND d.das <= 20 ) THEN 'Stage 2' "
471                      + "WHEN ( d.das > 20 AND d.das <= 30 ) THEN 'Stage 3' "
472                      + "WHEN ( d.das > 30 AND d.das <= 40 ) THEN 'Stage 4' "
473                      + "WHEN ( d.das > 40 AND d.das <= 50 ) THEN 'Stage 5' "
```

```
474                    + "ELSE 'Stage 6' END Stage "
475                    + "FROM imageev AS i, imgobjectev AS o, soyidentification
    AS s, dasplusev AS d "
476                    + "WHERE i.assayid = o.assayid "
477                    + "AND i.imgid = o.imgid "
478                    + "AND s.barcode = ( CAST( i.barcode AS INTEGER ) ) "
479                    + "AND i.assayid = d.assayid "
480                    + "AND i.fdate = d.fdate "
481                    + "AND i.set = d.set "
482                    + "AND ( s.line = 1 OR s.line = 2 OR s.line = 3 ) "
483                    + "AND i.camera = 'vis-side-1-0' "
484                    + "AND i.set = '2' "
485                    + "UNION "
486                    + "SELECT s.barcode, o.area, o.perimeter, o.circularity,
    "
487                    + "o.compactness, o.major, o.minor, o.eccentricity,
    o.hisgreypeak, "
488                    + "o.q1grey, o.q2grey, o.q3grey, "
489                    + "o.q1r, o.q2r, o.q3r,"
490                    + " o.q1g, o.q2g, o.q3g, "
491                    + "o.q1b, o.q2b, o.q3b, d.das, "
492                    + "CASE WHEN ( d.das <= 10 ) THEN 'Stage 1' "
493                    + "WHEN ( d.das > 10 AND d.das <= 20 ) THEN 'Stage 2' "
494                    + "WHEN ( d.das > 20 AND d.das <= 30 ) THEN 'Stage 3' "
495                    + "WHEN ( d.das > 30 AND d.das <= 40 ) THEN 'Stage 4' "
496                    + "WHEN ( d.das > 40 AND d.das <= 50 ) THEN 'Stage 5' "
497                    + "ELSE 'Stage 6' END Stage "
498                    + "FROM imageev AS i, imgobjectev AS o, soyidentification
    AS s, dasplusev AS d "
499                    + "WHERE i.assayid = o.assayid "
500                    + "AND i.imgid = o.imgid "
501                    + "AND s.barcode = ( CAST( i.barcode AS INTEGER ) ) "
502                    + "AND i.assayid = d.assayid "
503                    + "AND i.fdate = d.fdate "
504                    + "AND i.set = d.set "
505                    + "AND (s.line = 2 OR s.line = 3 ) "
506                    + "AND i.camera = 'vis-side-1-0' "
507                    + "AND i.set = '3' ";
508
509            PreparedStatement trainingPs = conn.prepareStatement
    (trainingSql);
510            ResultSet trainingSet = trainingPs.executeQuery();
511            while(trainingSet.next()) {
512                Double barcode = trainingSet.getDouble("barcode");
513                Double area = trainingSet.getDouble("area");
514                Double perimeter = trainingSet.getDouble("perimeter");
515                Double circularity = trainingSet.getDouble("circularity");
516                Double compactness = trainingSet.getDouble("compactness");
```

```
517                  Double major = trainingSet.getDouble("major");
518                  Double minor = trainingSet.getDouble("minor");
519                  Double eccentricity = trainingSet.getDouble("eccentricity");
520                  Double hisgreypeak = trainingSet.getDouble("hisgreypeak");
521                  Double q1grey = trainingSet.getDouble("q1grey");
522                  Double q2grey = trainingSet.getDouble("q2grey");
523                  Double q3grey = trainingSet.getDouble("q3grey");
524                  Double q1r = trainingSet.getDouble("q1r");
525                  Double q2r = trainingSet.getDouble("q2r");
526                  Double q3r = trainingSet.getDouble("q3r");
527                  Double q1g = trainingSet.getDouble("q1g");
528                  Double q2g = trainingSet.getDouble("q2g");
529                  Double q3g = trainingSet.getDouble("q3g");
530                  Double q1b = trainingSet.getDouble("q1b");
531                  Double q2b = trainingSet.getDouble("q2b");
532                  Double q3b = trainingSet.getDouble("q3b");
533                  Double das = trainingSet.getDouble("das");
534                  String stage = trainingSet.getString("Stage");
535
536                  //writing all the different attributes from the query to the
     output arff file
537                  trainingPw.write(barcode + "," + area + ", " + perimeter + ",
     " + circularity + ", " + compactness + ", " + major + ", " + minor + ", " +
     eccentricity
538                         + ", " + hisgreypeak + ", " + q1grey + ", " + q2grey
     + ", " + q3grey + ", " + q1r + ", " + q2r + ", " + q3r
539                         + ", " + q1g + ", " + q2g + ", " + q3g + ", " + q1b +
     ", " + q2b + ", " + q3b + "," + das + "," + "'" + stage + "'" + "\n");
540              }
541          trainingSet.close();
542          trainingPw.close();
543
544          PreparedStatement  testingPs = conn.prepareStatement(testingSql);
545          ResultSet testingSet = testingPs.executeQuery();
546          while(testingSet.next()) {
547              Double barcode = testingSet.getDouble("barcode");
548              Double area = testingSet.getDouble("area");
549              Double perimeter = testingSet.getDouble("perimeter");
550              Double circularity = testingSet.getDouble("circularity");
551              Double compactness = testingSet.getDouble("compactness");
552              Double major = testingSet.getDouble("major");
553              Double minor = testingSet.getDouble("minor");
554              Double eccentricity = testingSet.getDouble("eccentricity");
555              Double hisgreypeak = testingSet.getDouble("hisgreypeak");
556              Double q1grey = testingSet.getDouble("q1grey");
557              Double q2grey = testingSet.getDouble("q2grey");
558              Double q3grey = testingSet.getDouble("q3grey");
559              Double q1r = testingSet.getDouble("q1r");
```

```java
560                 Double q2r = testingSet.getDouble("q2r");
561                 Double q3r = testingSet.getDouble("q3r");
562                 Double q1g = testingSet.getDouble("q1g");
563                 Double q2g = testingSet.getDouble("q2g");
564                 Double q3g = testingSet.getDouble("q3g");
565                 Double q1b = testingSet.getDouble("q1b");
566                 Double q2b = testingSet.getDouble("q2b");
567                 Double q3b = testingSet.getDouble("q3b");
568                 Double das = testingSet.getDouble("das");
569                 String stage = testingSet.getString("Stage");
570
571                 testingPw.write(barcode + "," + area + ", " + perimeter + ",
    " + circularity + ", " + compactness + ", " + major + ", " + minor + ", " +
    eccentricity
572                             + ", " + hisgreypeak + ", " + q1grey + ", " + q2grey
    + ", " + q3grey + ", " + q1r + ", " + q2r + ", " + q3r
573                             + ", " + q1g + ", " + q2g + ", " + q3g + ", " + q1b +
    ", " + q2b + ", " + q3b + "," + das + "," + "'" + stage + "'" + "\n");
574             }
575             testingSet.close();
576             testingPw.close();
577
578             conn.close();
579         }
580         catch (ClassNotFoundException e) {
581
582             System.out.println("Improper database connection set-up.");
583             e.printStackTrace();
584
585         }
586     }
587
588 public static Instances [] attributeSelector(Instances train, Instances
    test) throws Exception {
589
590         AttributeSelection selector = new AttributeSelection();
591         Ranker search = new Ranker();
592         OneRAttributeEval eval = new OneRAttributeEval();
593         selector.setEvaluator(eval);
594         selector.setSearch(search);
595         selector.SelectAttributes(train);
596
597         //rankedAttributes gives the ranking of the attributes along with
    their weights
598         //selected Attributes just gives the order of the ranking
599         Instances trainTemp = selector.reduceDimensionality(train);
600         Instances trainTest = selector.reduceDimensionality(test);
601
```

```
602
603         Instances [] reduced = {trainTemp, trainTest};
604
605         return reduced;
606    }
607
608 }
```

```r
library(ggplot2)
library(reshape2)
setwd("~/Documents/U4/Comp401/output")

predictions <- read.csv("originalPredictions.csv", header = TRUE)
unswitched <- rep(0, length(predictions) - 3)
classifier.name <- colnames(predictions[4:length(predictions)])
actual.stage <- as.vector(predictions[, "Actual"])
n <- length(actual.stage)

for(i in 4:length(predictions)){
    v <- as.vector(predictions[,i])
    unswitched[i-3] <- round(100*length(which(actual.stage == v))/n, 3)
}


### RESULTS WHEN TRAINING DATA AND TESTING DATA ARE SWITCHED
switchedPredictions <- read.csv("switchedPredictions.csv", header = TRUE)
switched <- rep(0, length(switchedPredictions) - 3)
classifier.name.s <- colnames(switchedPredictions[4:length(switchedPredictions)])
actual.stage.s <- as.vector(switchedPredictions[, "Actual"])
sn <- length(actual.stage.s)

for(i in 4:length(switchedPredictions)){
  sv <- as.vector(switchedPredictions[,i])
  switched[i-3] <- round(100*length(which(actual.stage.s == sv))/sn, 3)
}

allPrecisions <- data.frame(classifier.name, unswitched, switched)
allPrecisions<- melt(allPrecisions, id.vars = "classifier.name")
precision.plot <- ggplot(allPrecisions, aes(y = value, x = classifier.name, fill = variable))
precision.plot <- precision.plot + geom_bar(stat = "identity", position = "dodge")
precision.plot <- precision.plot + xlab("Classifier Name") + ylab("Percent Precision")
precision.plot <- precision.plot + guides(fill=guide_legend(title = "Training Data"))
precision.plot <- precision.plot + scale_fill_manual(values = c("#0033FF", "#6699FF"))
precision.plot <- precision.plot + geom_text(aes(label = round(value, 1)), position = position_dodge(0.75), vjust = -1)
precision.plot <- precision.plot + theme(axis.text.x = element_text(angle = 45, hjust = 1), text = element_text(size = 24))
precision.plot

### Finding out which algorithms performed better between switched data
###BETTER UNSWITCHED
classifier.name[which(unswitched > switched)]

###BETTER SWITCHED
classifier.name[which(unswitched < switched)]

### RESULTS WHEN ATTRIBUTES ARE RANKED FIRST
rankedPredictions <- read.csv("RankedPredictions.csv", header = TRUE)
ranked <- rep(0, length(rankedPredictions) - 3)
classifier.name.r <- colnames(rankedPredictions[4:length(rankedPredictions)])
actual.stage.r <- as.vector(rankedPredictions[, "Actual"])
rn <- length(actual.stage.r)

for(i in 4:length(rankedPredictions)){
  rv <- as.vector(rankedPredictions[,i])
  ranked[i-3] <- round(100*length(which(actual.stage.r == rv))/rn, 3)
}

allPrecisions <- data.frame(classifier.name.r, switched, ranked)
allPrecisions<- melt(allPrecisions, id.vars = "classifier.name.r")
precision.plot <- ggplot(allPrecisions, aes(y = value, x = classifier.name.r, fill = variable))
precision.plot <- precision.plot + geom_bar(stat = "identity", position = "dodge")
precision.plot <- precision.plot + xlab("Classifier Name") + ylab("Percent Precision")
precision.plot <- precision.plot + guides(fill=guide_legend(title = "Training Data"))
precision.plot <- precision.plot + scale_fill_manual(values = c("#0033FF", "#6699FF"))
precision.plot <- precision.plot + geom_text(aes(label = round(value, 1)), position = position_dodge(0.75), vjust = -1)
precision.plot <- precision.plot + theme(axis.text.x = element_text(angle = 45, hjust = 1), text = element_text(size = 24))
precision.plot

###BETTER RANKED
classifier.name.r[which(ranked > switched)]

###NO DIFFERENCE
classifier.name.r[which(ranked == switched)]

###BETTER UNRANKED
classifier.name.r[which(ranked < switched)]

###LOOKING AT THE STAGE PREDICTION BREAKDOWN WHEN THE ATTRIBUTES ARE RANKED
predictions <- read.csv("RankedPredictions.csv", header = TRUE)
classifier.name <- colnames(predictions[5:length(predictions)])
actual.stage <- as.vector(predictions[, "Actual"])
n <- length(actual.stage)

stages <- c("Stage 1", "Stage 2", "Stage 3", "Stage 4", "Stage 5")
numStages <- rep(0, length(stages))
for (i in 1:length(stages)){
  numStages[i] <- length(which(predictions[,"Actual"] == stages[i]))
```

```r
}

percent.correct <- rep(0, length(stages))

for(i in 1:length(stages)){
  indices <- which(predictions[,"Actual"] == stages[i])
  for(j in 5:length(predictions)){
    percent.correct[i] <- percent.correct[i] + length(which(predictions[indices, j] == stages[i]))
  }
}

### percent.correct is an AVERAGE of the percent time a stage is predicted correctly
percent.correct <- 100*(round(percent.correct/(numStages*(length(classifier.name)-1)),3))
percent.correct <- data.frame(stages, percent.correct)

stage.plot <- ggplot(percent.correct, aes(y = percent.correct, x = stages))
stage.plot <- stage.plot + geom_bar(stat = "identity", fill = "#0033FF")
stage.plot <- stage.plot + xlab("Stage") + ylab("Percent Correctly Predicted") + ggtitle("Average Precision per Stage")
stage.plot <- stage.plot + geom_text(aes(label = round(percent.correct, 1)), position = position_dodge(0.75), vjust =
-1)
stage.plot

###What percentage of the time was stage 1 predicted?
times.predicted <- rep(0, length(stages))

for(i in 1:length(stages)){
  for(j in 5:length(predictions)){
    times.predicted[i] <- times.predicted[i] + length(which(predictions[,j] == stages[i]))
  }
}

times.predicted <- 100*(round(times.predicted/(sum(times.predicted)),5))
times.predicted <- data.frame(stages, times.predicted)
times.predicted.plot <- ggplot(times.predicted, aes(y = times.predicted, x = stages))
times.predicted.plot <- times.predicted.plot + geom_bar(stat = "identity", fill = "#0033FF")
times.predicted.plot <- times.predicted.plot + xlab("Stage") + ylab("Percent of Time Predicted") + ggtitle("Percentage
Predicted per Stage")
times.predicted.plot <- times.predicted.plot + geom_text(aes(label = round(times.predicted, 1)), position =
position_dodge(0.75), vjust = -1)
times.predicted.plot

### WORKING WITH JUST THE SWITCHED DATA SINCE THIS LED TO OVERALL BETTER RESULTS
predictions <- read.csv("switchedPredictions.csv", header = TRUE)
classifier.name <- colnames(predictions[5:length(predictions)])
actual.stage <- as.vector(predictions[, "Actual"])
n <- length(actual.stage)

stages <- c("Stage 1", "Stage 2", "Stage 3", "Stage 4", "Stage 5")
numStages <- rep(0, length(stages))
for (i in 1:length(stages)){
  numStages[i] <- length(which(predictions[,"Actual"] == stages[i]))
}

percent.correct <- rep(0, length(stages))

for(i in 1:length(stages)){
  indices <- which(predictions[,"Actual"] == stages[i])
  for(j in 5:length(predictions)){
    percent.correct[i] <- percent.correct[i] + length(which(predictions[indices, j] == stages[i]))
  }
}

### percent.correct is an AVERAGE of the percent time a stage is predicted correctly
percent.correct <- 100*(round(percent.correct/(numStages*(length(classifier.name)-1)),3))
percent.correct <- data.frame(stages, percent.correct)

stage.plot <- ggplot(percent.correct, aes(y = percent.correct, x = stages))
stage.plot <- stage.plot + geom_bar(stat = "identity", fill = "#0033FF")
stage.plot <- stage.plot + xlab("Stage") + ylab("Percent Correctly Predicted")
stage.plot <- stage.plot + geom_text(aes(label = round(percent.correct, 1)), position = position_dodge(0.75), vjust =
-1)
stage.plot <- stage.plot + theme(axis.text.x = element_text(angle = 45, hjust = 1), text = element_text(size = 18))
stage.plot

###What percentage of the time was stage 1 predicted?
predictions <- read.csv("switchedPredictions.csv", header = TRUE)
times.predicted <- rep(0, length(stages))

for(i in 1:length(stages)){
  for(j in 5:length(predictions)){
    times.predicted[i] <- times.predicted[i] + length(which(predictions[,j] == stages[i]))
  }
}

times.predicted <- 100*(round(times.predicted/(sum(times.predicted)),5))
times.predicted <- data.frame(stages, times.predicted)
times.predicted.plot <- ggplot(times.predicted, aes(y = times.predicted, x = stages))
times.predicted.plot <- times.predicted.plot + geom_bar(stat = "identity", fill = "#0033FF")
times.predicted.plot <- times.predicted.plot + xlab("Stage") + ylab("Percent of Time Predicted")
times.predicted.plot <- times.predicted.plot + geom_text(aes(label = round(times.predicted, 1)), position =
position_dodge(0.75), vjust = -1)
```

```r
times.predicted.plot <- times.predicted.plot + theme(axis.text.x = element_text(angle = 45, hjust = 1), text =
element_text(size = 18))
times.predicted.plot

###How many of each stage was there in the testing set?
trainingSet <- read.csv("DatabaseTesting.csv", header = TRUE)
stages.training <- rep(0, length(stages))
for (i in 1:length(stages)){
  stages.training[i] <- length(which(trainingSet[,"Stage"] == stages[i]))
}


####WORKING WITH THE ORIGINAL DATA, COPIED CODE BASICALLY FROM ABOVE
predictions <- read.csv("originalPredictions.csv", header = TRUE)
classifier.name <- colnames(predictions[5:length(predictions)])
actual.stage <- as.vector(predictions[, "Actual"])
n <- length(actual.stage)

stages <- c("Stage 1", "Stage 2", "Stage 3", "Stage 4", "Stage 5")
numStages <- rep(0, length(stages))
for (i in 1:length(stages)){
  numStages[i] <- length(which(predictions[,"Actual"] == stages[i]))
}

percent.correct <- rep(0, length(stages))

for(i in 1:length(stages)){
  indices <- which(predictions[,"Actual"] == stages[i])
  for(j in 5:length(predictions)){
    percent.correct[i] <- percent.correct[i] + length(which(predictions[indices, j] == stages[i]))
  }
}

### percent.correct is an AVERAGE of the percent time a stage is predicted correctly
percent.correct <- 100*(round(percent.correct/(numStages*(length(classifier.name)-1)),3))
percent.correct <- data.frame(stages, percent.correct)

stage.plot <- ggplot(percent.correct, aes(y = percent.correct, x = stages))
stage.plot <- stage.plot + geom_bar(stat = "identity", fill = "#0033FF")
stage.plot <- stage.plot + xlab("Stage") + ylab("Percent Correctly Predicted") + ggtitle("Average Precision per Stage")
stage.plot <- stage.plot + geom_text(aes(label = round(percent.correct, 1)), position = position_dodge(0.75), vjust =
-1)
stage.plot

###What percentage of the time was stage 1 predicted?
times.predicted <- rep(0, length(stages))

for(i in 1:length(stages)){
  for(j in 5:length(predictions)){
    times.predicted[i] <- times.predicted[i] + length(which(predictions[,j] == stages[i]))
  }
}

times.predicted <- 100*(round(times.predicted/(sum(times.predicted)),5))
times.predicted <- data.frame(stages, times.predicted)
times.predicted.plot <- ggplot(times.predicted, aes(y = times.predicted, x = stages))
times.predicted.plot <- times.predicted.plot + geom_bar(stat = "identity", fill = "#0033FF")
times.predicted.plot <- times.predicted.plot + xlab("Stage") + ylab("Percent of Time Predicted") + ggtitle("Percentage
Predicted per Stage")
times.predicted.plot <- times.predicted.plot + geom_text(aes(label = round(times.predicted, 1)), position =
position_dodge(0.75), vjust = -1)
times.predicted.plot
```