# CSED/AIGS 526: Data Mining

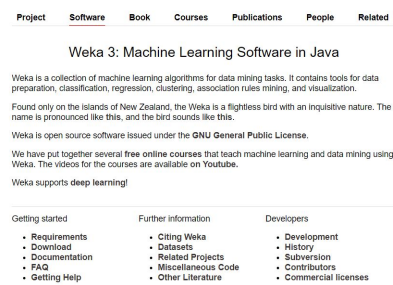**Your name:** _____ Jaeyoon Sim _____

**Student ID:** _____ 20222421 _____

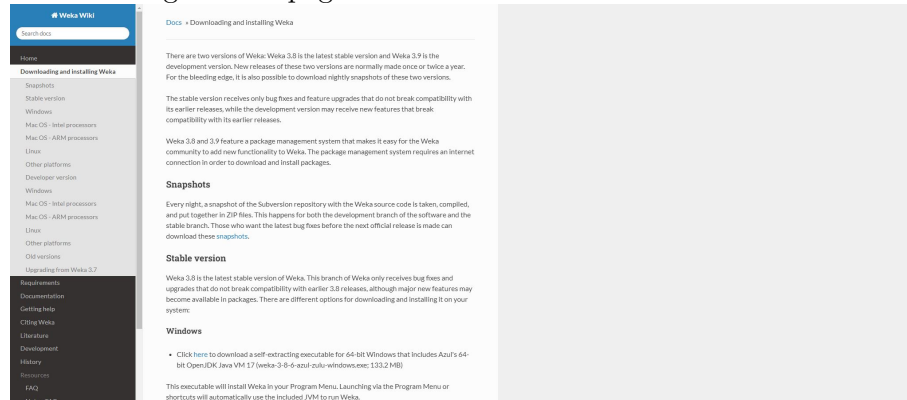**Use late days quota(Yes/No):** _____ No _____

# Problem 1

## Weka

The purpose of this assignment is to learn about a program called "Weka" and to get used to the data given. "Weka" is a collection of machine learning algorithms for data mining tasks. It contains tools for data preparation, classification, regression, clustering, association rules mining, and visualization. That's why we should install "Weka" on our computer.

1. We access https://www.cs.waikato.ac.nz/ml/weka/.
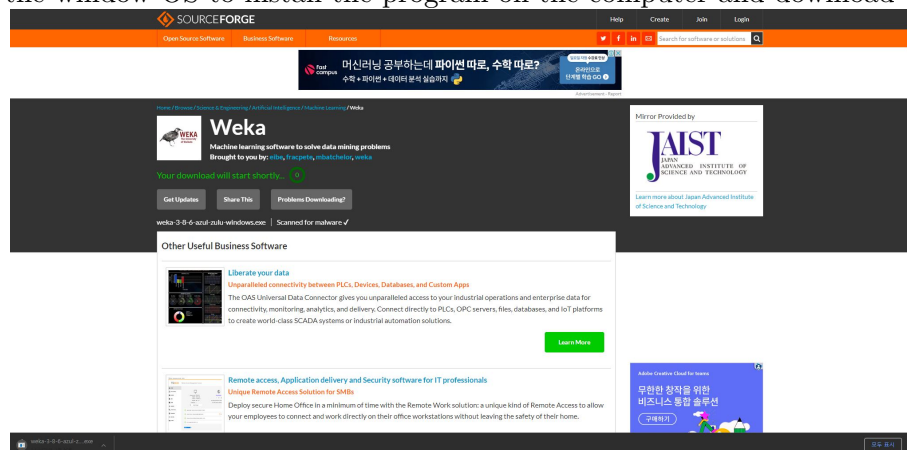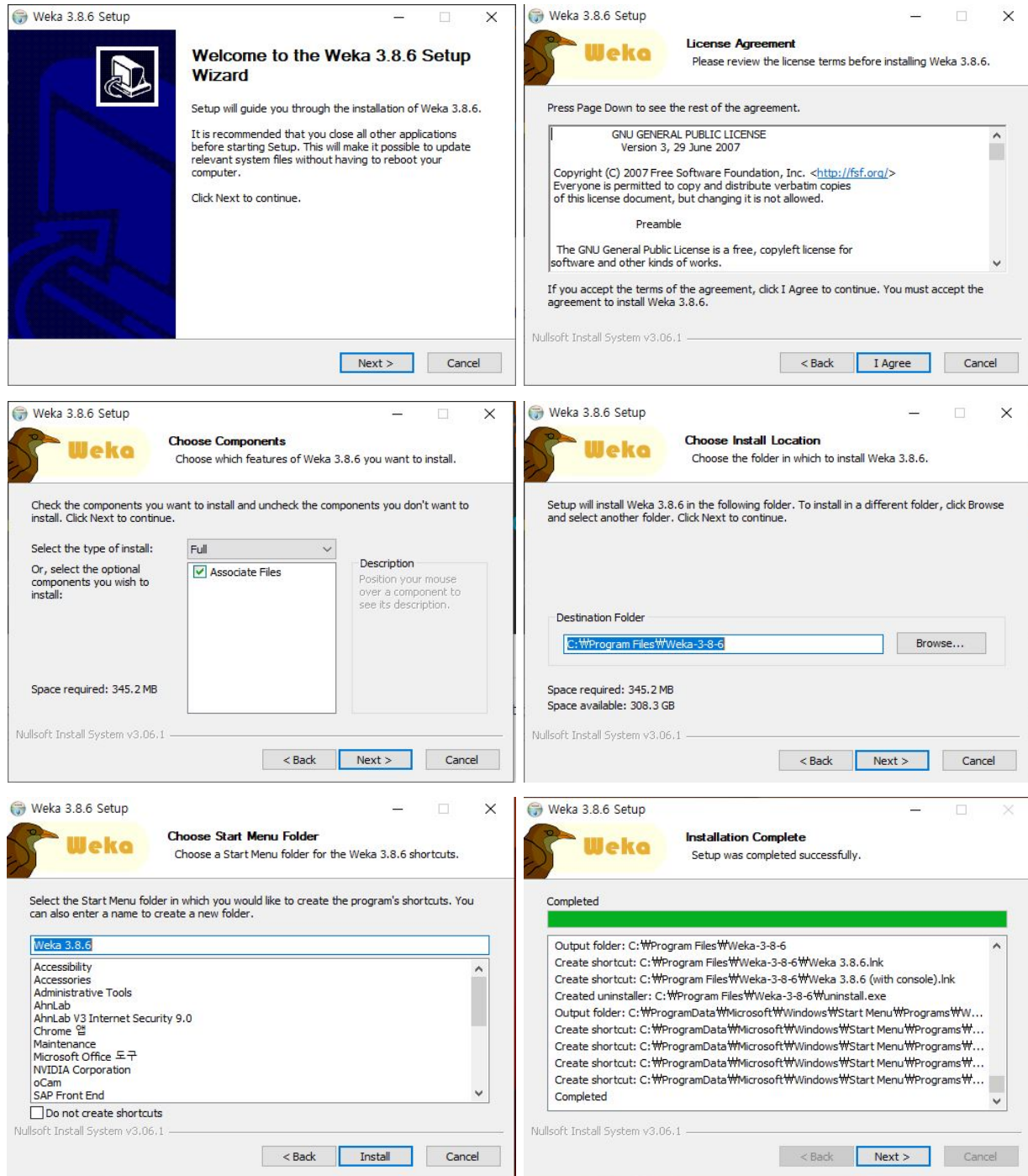


2. We click Download to go to the page where the "Weka" document is located.



3. We select the window OS to install the program on the computer and download the program.

4. From now on, we use the .exe file to install "Weka".

5. When we complete the installation of the "Weka" program, we can run the "Weka" program as follows.

## Problem 2

### Preprocess, and Clustering

We received an .arff (attribute-relation file format) file used by "Weka" for this assignment. In order to use this .arff file, we open "Weka" and click "Explorer" applications. Data files for "Weka" is .arff, and we can open this file in "Weka" by following the below.

With "nbastats2021.arff" file, the dataset contains basketball related records from 240 players from National Basketball League(NBA). We will look at the data set made based on the various information commonly used in basketball games. Before that, we would like to briefly look at the attributes used in the "Weka" program.

1. GP: Game Played

2. MPG: Minutes played per game

3. FT%: Free throw percentage

4. 2P%: 2-points field goal percentage

5. 3P%: 3-points field goal percentage

6. PPG: Points per game

7. RPG: Rebounds per game

8. APG: Assists per game

9. SPG: Steals per game

10. BPG: Blocks per game

11. TOPG: Turnovers per game

12. POS: Position

Given nbastat2021 dataset contains statistics from National Basketball League(NBA). Each player is characterized by various features seen above. In particular, attribute 1 to 11 are features of players related to basketball games, and feature 12 is about positions in which players are played in basketball games.

In actual basketball games, roles are divided from 1 to 5 depending on the position. Number one is a point guard, number two is a shooting guard, number three is a small forward, number four is a power forward, and number five is a center. Each role has different things to do, and we would like to use positions such as {G, G-F, F, F-C, C} in this assignment.

First of all, we want to look at the data in the "Preprocess" tab of the "Weka" program. The "Preprocess" tab currently contains information about the dataset. We can find out what data sets are, how many attributes are in total, and how many instances are in total. The following is the result of visualization by creating a histogram from all attributes of the data.

Of the total 12 histograms, we will look closely at the histograms of MPG. If we look at the MPG histogram, we can see how long 240 players played the game.



It expresses 240 players in histograms every 2.6 minutes from a minimum of 20 minutes to a maximum of 38.2 minutes. The players played basketball for an average of 28.464 minutes, with a standard deviation of 4.805 minutes.

Now, we want to look at the data in the "Cluster" tab of the "Weka" program. Considering the results of visualizing histograms from all attributes of data and the roles given for each position in basketball games, we would like to proceed with clustering using {3P%, RPG, APG, BPG}. And we would like to exclude {GP, MPG, FT%, 2P%, PPG, SPG, TOP, POS}. The reason is that once we think about our personal experience and role by position in basketball games, rebounds and blocks match the position of the center, while 3P% and assist match the position of the guard. In general, this trend is shown, but considering the given data, it tends to be consistent to some extent, and overall, there is no significant difference in histograms for each position in the number of games played, time of game, free throw, 2-points field goal, scoring, still, and turnover.

This time, we're going to use the "SimpleKMeans" clustering method stored in the "Weka" Cluster tab, ignoring the attributes we mentioned earlier. In the case of the KMean algorithm, it is sensitive to the distribution of data and has weak adaptability to the Outlier. The "SimpleKMeans" algorithm in "Weka" uses Euclidean distance measure to compute distances between instances and clusters. Here, the "POS" attribute will be excluded when attempting to perform the clustering method. Along with these conditions, the total number of clusters is set to five to cluster the positions of basketball players.

The following result is the result of excluding only "POS" and executing from the attributes. As a result of SimpleKMeans, it can be seen that 166 out of 240 instances are not accurately clustered. So it can be seen that 69.1667% of the total instances did not properly form a cluster excluding only "POS".

This is the result of clustering using {3P%, RPG, APG, BPG} attributes as we initially thought. This time, as a result, it can be seen that 135 out of 240 instances are not accurately clustered. Therefore, 56.25% of the total instances did not properly form clusters. It can be seen that this resulted in better clustering results based on existing thoughts and analysis.

The final problem of Problem 2 is to create five clusters again and visualize them using the "EM" clustering algorithm. In statistics, an expectation–maximization (EM) algorithm is an iterative method to find (local) maximum likelihood or maximum a posteriori(MAP) estimates of parameters in statistical models, where the model depends on unobserved latent variables. The EM iteration alternates between performing an expectation(E) step, which creates a function for the expectation of the log-likelihood evaluated using the current estimate for the parameters, and a maximization(M) step, which computes parameters maximizing the expected log-likelihood found on the E step. These parameter-estimates are then used to determine the distribution of the latent variables in the next E step.

In this problem, using the EM algorithm, we can also change the hyperparameters of EM cluster. First of all, we want to set the EM clustering method to different 'ignored attributes' and execute it. EM clustering was performed by including {3P%, RPG, APG, BPG} attributes in the first case, and EM clustering was performed by including {GP, MPG, FT%, 2P%, PPG, SPG, TOP} attributes. In the first case, {GP, MPG, FT%, 2P%, PPG, SPG, TOP, POS} attributes were ignored, and in the second case, {3P%, RPG, APG, BPG, POS} attributes were ignored.

Clusterer output

=== Run information ===

Scheme:    weka.clusterers.EM -I 100 -N 5 -X 10 -max -1 -ll-cv 1.0E-6 -ll-iter 1.0E-6 -M 1.0E-6 -K 10 -n
Relation:  nbastat2021
Instances: 240
Attributes: 12
           GP
           MPG
           FT%
           2P%
           PPG
           SPG
           TOP
Ignored:
           3P%
           RPG
           APG
           BPG
           POS
Test mode: Classes to clusters evaluation on training data

=== Clustering model (full training set) ===

EM
==

Number of clusters: 5
Number of iterations performed: 40

               Cluster
Attribute    0     1     2     3     4
           (0.23) (0.25) (0.22) (0.09) (0.22)
==================================================
GP
  mean     50.3719 61.705 54.6967 43.1888 57.0739
  std. dev. 11.6361 5.1365 10.1218 11.5963 9.6168

MPG
  mean     22.4925 26.543 34.526 27.3495 31.0517
  std. dev. 1.4958 2.527 1.4766 2.6621 2.1592

FT%
  mean     0.7416 0.7533 0.8127 0.8001 0.7997
  std. dev. 0.0882 0.1156 0.0611 0.0899 0.0766

2P%
  mean     0.5535 0.5399 0.5182 0.4834 0.5291
  std. dev. 0.0741 0.0758 0.0472 0.046 0.0629

---

Clusterer output

           std. dev. 0.0741 0.0758 0.0472 0.046 0.0629

PPG
  mean     8.326 10.219 21.7289 11.9066 15.8968
  std. dev. 1.8022 2.0198 4.5322 2.9275 1.8211

SPG
  mean     0.6511 0.8456 1.2107 1.1006 0.8032
  std. dev. 0.2461 0.3029 0.3733 0.3563 0.1838

TOP
  mean     0.8372 1.0619 2.8144 1.8713 1.6796
  std. dev. 0.1937 0.3027 0.7607 0.7053 0.4111

Time taken to build model (full training data) : 0.02 seconds

=== Model and evaluation on training set ===

Clustered Instances

0    55 ( 23%)
1    60 ( 25%)
2    52 ( 22%)
3    17 (  7%)
4    56 ( 23%)

Log likelihood: -7.47534

Class attribute: POS
Classes to Clusters:

 0  1  2  3  4  <-- assigned to cluster
14 20 27 12 24 | G
14 11  7  3  5 | G-F
16 18 11  2 14 | F
11  6  6  0  6 | F-C
 0  5  1  0  7 | C

Cluster 0 <-- G-F
Cluster 1 <-- F
Cluster 2 <-- G
Cluster 3 <-- No class
Cluster 4 <-- C

Incorrectly clustered instances :      174.0     72.5    %

---

Weka Clusterer Visualize: 15:55:06 - EM (nbastat2021)

Now, we will try to make five clusters each representing a different position-{G, G-F, F, F-C, C} and visualize one of them. Like the SimpleKMeans above, clustering was performed with {3P%, RPG, APG, BPG} attributes left. The results were compared while adjusting jitter, which is the degree to which data is scattered. In the case of hyperparameter, it was used as default except for changing the number of clusters to five.

---

Clusterer output

           (0.1)   (0.25)  (0.26)  (0.29)  (0.1)
========================================
3P%
  mean      0.1642  0.3578  0.3598  0.3447  0.308
  std. dev. 0.1498  0.0405  0.0375  0.0428  0.1879

RPG
  mean      7.9634  4.5943  2.8189  5.0973  9.5307
  std. dev. 2.3198  1.0044  0.5761  1.5385  1.8721

APG
  mean      1.3628  1.7476  2.4169  5.214   3.9033
  std. dev. 0.5976  0.6329  1.0873  2.0191  1.7636

BPG
  mean      1.4296  0.4999  0.2253  0.4306  0.935
  std. dev. 0.5889  0.2456  0.1026  0.2022  0.3127

Time taken to build model (full training data) : 0.01 seconds

=== Model and evaluation on training set ===

Clustered Instances

0    25 ( 10%)
1    64 ( 27%)
2    65 ( 27%)
3    64 ( 27%)
4    22 (  9%)

Log likelihood: -2.60723

Class attribute: POS
Classes to Clusters:

 0  1  2  3  4  <-- assigned to cluster
 0 13 39 45  0 | G
 1 12 19  8  0 | G-F
 4 32  7 10  8 | F
13  7  0  1  8 | F-C
 7  0  0  0  6 | C

Cluster 0 <-- F-C
Cluster 1 <-- F
Cluster 2 <-- G-F
Cluster 3 <-- G
Cluster 4 <-- C

Incorrectly clustered instances :      125.0     52.0833 %

---

Weka Clusterer Visualize: 16:05:55 - EM (nbastat2021)

The results of clustering are not so satisfactory. The best result is that there are only instances corresponding to the position for each cluster, but our result was not. However, unlike SimpleKMeans clustering, which clustered 135 instances incorrectly, EM clustering, which clustered 125 instances incorrectly, showed better results.

# Problem 3

## Classification

As mentioned in Problem 1, each player is characterized by 11 features and a position. Now, we can classify the players by their positions using the given features in "Weka". And the "Weka" program provides several decision tree-based algorithm. In this problem, we will classify the data using C4.5 decision tree(J48 in "Weka") with default parameters. With using parameters, we will train the tree on the data and test on the same data.

The main feature of the C4.5 decision tree is that it first accepts both continuous and discrete variables. It can handle incomplete data, and solve the overfitting problem through the pruning process of cutting off the end of the decision tree. In addition, different weights may be applied to the variables constituting the training data. We will classify the data through this C4.5 decision tree algorithm. The following is the screenshot of summary and detailed accuarcy by class.

```
=== Summary ===

Correctly Classified Instances          207               86.25   %
Incorrectly Classified Instances         33               13.75   %
Kappa statistic                           0.8098
Mean absolute error                       0.0875
Root mean squared error                   0.2091
Relative absolute error                  30.0218 %
Root relative squared error              54.8473 %
Total Number of Instances               240

=== Detailed Accuracy By Class ===

                TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
                0.938    0.077    0.892      0.938   0.915      0.855  0.966     0.917     G
                0.725    0.040    0.784      0.725   0.753      0.707  0.931     0.764     G-F
                0.820    0.045    0.862      0.820   0.840      0.788  0.955     0.885     F
                0.897    0.024    0.839      0.897   0.867      0.848  0.988     0.877     F-C
                0.846    0.004    0.917      0.846   0.880      0.874  0.998     0.944     C
Weighted Avg.   0.863    0.052    0.861      0.863   0.861      0.814  0.962     0.880
```
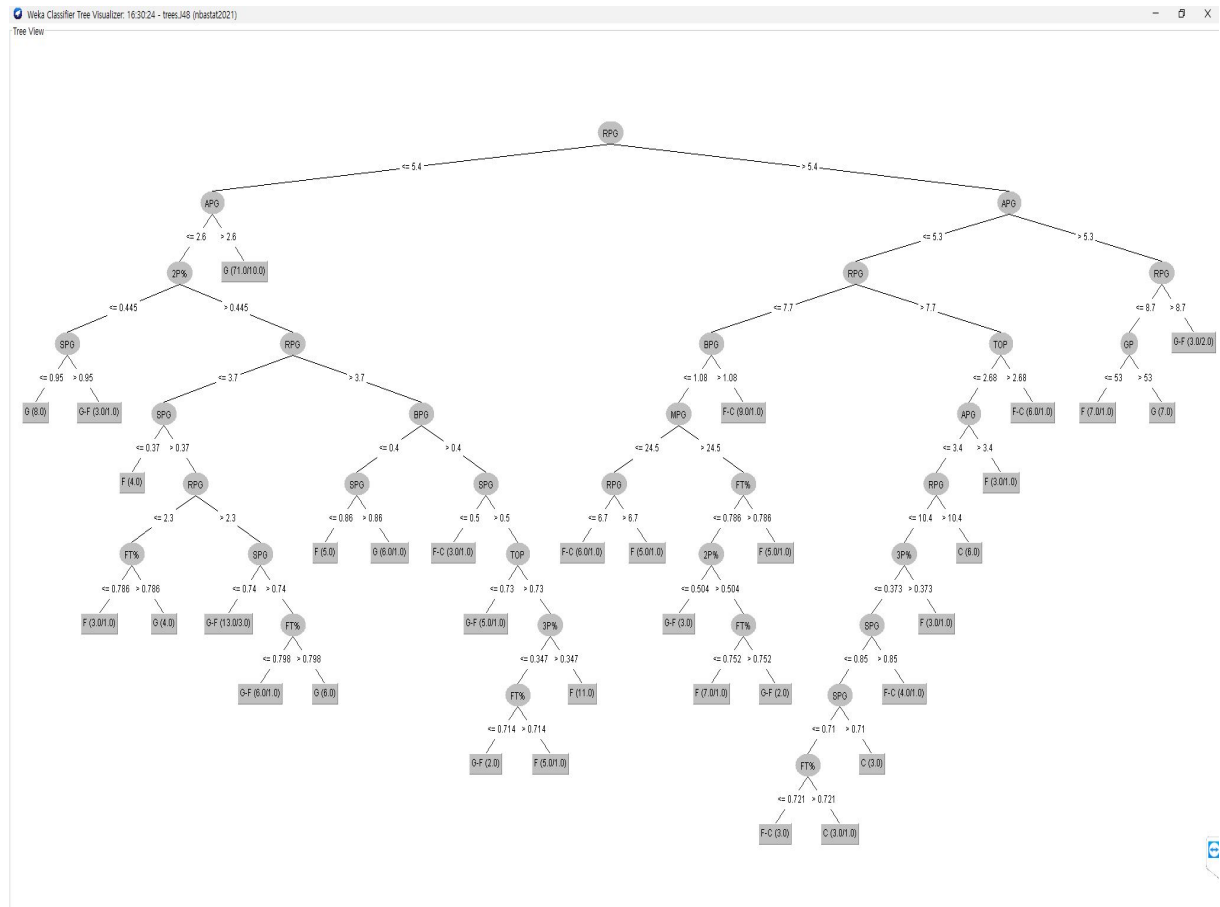
The following is a visualization of the results using decision tree. The size of the tree is the total number of nodes in the tree, which is terminal nodes + nonterminal nodes. In this case, the size of tree is 67.



```
=== Confusion Matrix ===

  a  b  c  d  e   <-- classified as
 91  3  3  0  0 |  a = G
  8 29  2  1  0 |  b = G-F
  3  4 50  4  0 |  c = F
  0  0  2 26  1 |  d = F-C
  0  1  1  0 11 |  e = C
```

The above is the confusion matrix provided by "Weka", where row shows the true class and column shows the predicated class. Of the instances with a total of 97 guard(G) positions, 91 were correctly classified. And 3 of the instances with a total of 61 forward(F) positions were incorrectly classified as guard(G).

Finally, in this problem, we can change the parameters of the C4.5 decision tree. We will classify using C4.5 decision tree with 50% pruning rate. The following is the screenshot of summary and detailed accuarcy by class.

```
=== Summary ===

Correctly Classified Instances         212               88.3333 %
Incorrectly Classified Instances        28               11.6667 %
Kappa statistic                          0.8394
Mean absolute error                      0.0724
Root mean squared error                  0.1902
Relative absolute error                 24.8399 %
Root relative squared error             49.8898 %
Total Number of Instances              240


=== Detailed Accuracy By Class ===

                TP Rate  FP Rate  Precision  Recall  F-Measure  MCC     ROC Area  PRC Area  Class
                0.938    0.042    0.938      0.938   0.938      0.896   0.984     0.964     G
                0.775    0.040    0.795      0.775   0.785      0.743   0.955     0.818     G-F
                0.869    0.045    0.869      0.869   0.869      0.824   0.976     0.920     F
                0.897    0.024    0.839      0.897   0.867      0.848   0.988     0.877     F-C
                0.846    0.004    0.917      0.846   0.880      0.874   0.998     0.944     C
Weighted Avg.   0.883    0.038    0.883      0.883   0.883      0.845   0.978     0.917
```
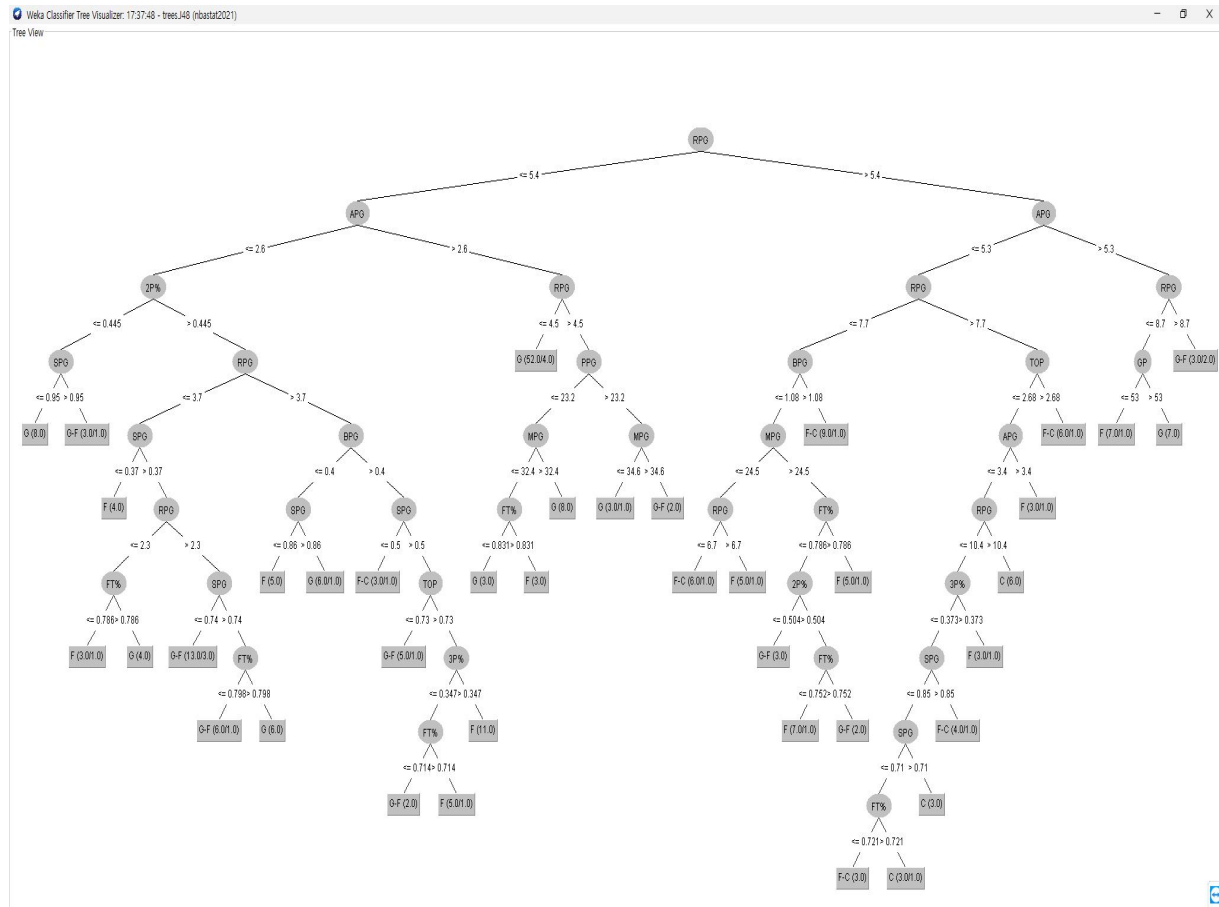
```
=== Confusion Matrix ===

  a  b  c  d  e    <-- classified as
 91  3  3  0  0 |  a = G
  6 31  2  1  0 |  b = G-F
  0  4 53  4  0 |  c = F
  0  0  2 26  1 |  d = F-C
  0  1  1  0 11 |  e = C
```

I think RPG is the most important attribute and FT% is the least important attribute. In the Decision tree, the overall tree shape can be complicated or simplified depending on the attribute selected first. As the tree shape becomes more complex, the attribute that can create a simple tree becomes more important because overfitting occurs more easily. We can also find informative attributes directly in the same way as information gain, gini impurity, etc. However, it was judged that the optimal result was generated when the decision tree was used with training data in "Weka", although it was omitted because there were many attributes and computations.