# University of Colombo School of Computing

# Bachelor of Information Systems

# Data Analytics

# IS 4003

# Assignment – Data Mining

E.H.Grero
15020258
2015/IS/025

# Objectives

The main objective of this dataset is to measure geometrical properties of kernels belonging to three different varieties of wheat namely Kama, Rosa and Canadian. A soft X-ray technique and GRAINS package is been used to construct all seven, real-valued attributes.

Benefit of using association rule mining would be to identify the inter relationships of the geometrical properties and identify separately the properties that will be able to identify the three different varieties of wheat Kama, Rosa and Canadian separately and clearly.

# Dataset description

The dataset was taken from the UCI machine learning repository.

URL of the dataset :- https://archive.ics.uci.edu/ml/datasets/seeds

The source of the dataset is taken from MaÅ, gorzata Charytanowicz, Jerzy Niewczas from Institute of Mathematics and Computer Science and Piotr Kulczycki, Piotr A. Kowalski, Szymon Lukasik, Slawomir Zak from Department of Automatic Control and Information Technology

To construct the data, seven geometric parameters of wheat kernels have been measured:

1. area A,
2. perimeter P,
3. compactness C = 4*pi*A/P^2,
4. length of kernel,
5. width of kernel,
6. asymmetry coefficient
7. length of kernel groove.

All these parameters are real-valued continuous.

The examined group or the data set comprised kernels belonging to three different varieties of wheat: Kama, Rosa and Canadian, 70 elements each, randomly selected for the experiment. High quality visualization of the internal kernel structure has been detected using a soft X-ray technique. It is non-destructive and considerably cheaper than other more sophisticated imaging techniques like scanning microscopy or laser technology. The images were recorded on 13x18 cm X-ray KODAK plates. Studies were conducted using combine harvested wheat grain originating from experimental fields, explored at the Institute of Agrophysics of the Polish Academy of Sciences in Lublin.

## Reason for selecting the dataset

The dataset contains 210 instances. 70 each from the three varieties of wheat namely Kama, Rosa and Canadian. Therefore, the dataset contains enough instances to identify association rules and patterns. The data set can be used for the tasks of classification and cluster analysis also made me choose this dataset.

## Preprocessing

Dataset when downloaded from the UCI machine repository was a text file. First the text file was open from Microsoft excel and it is shown as below.
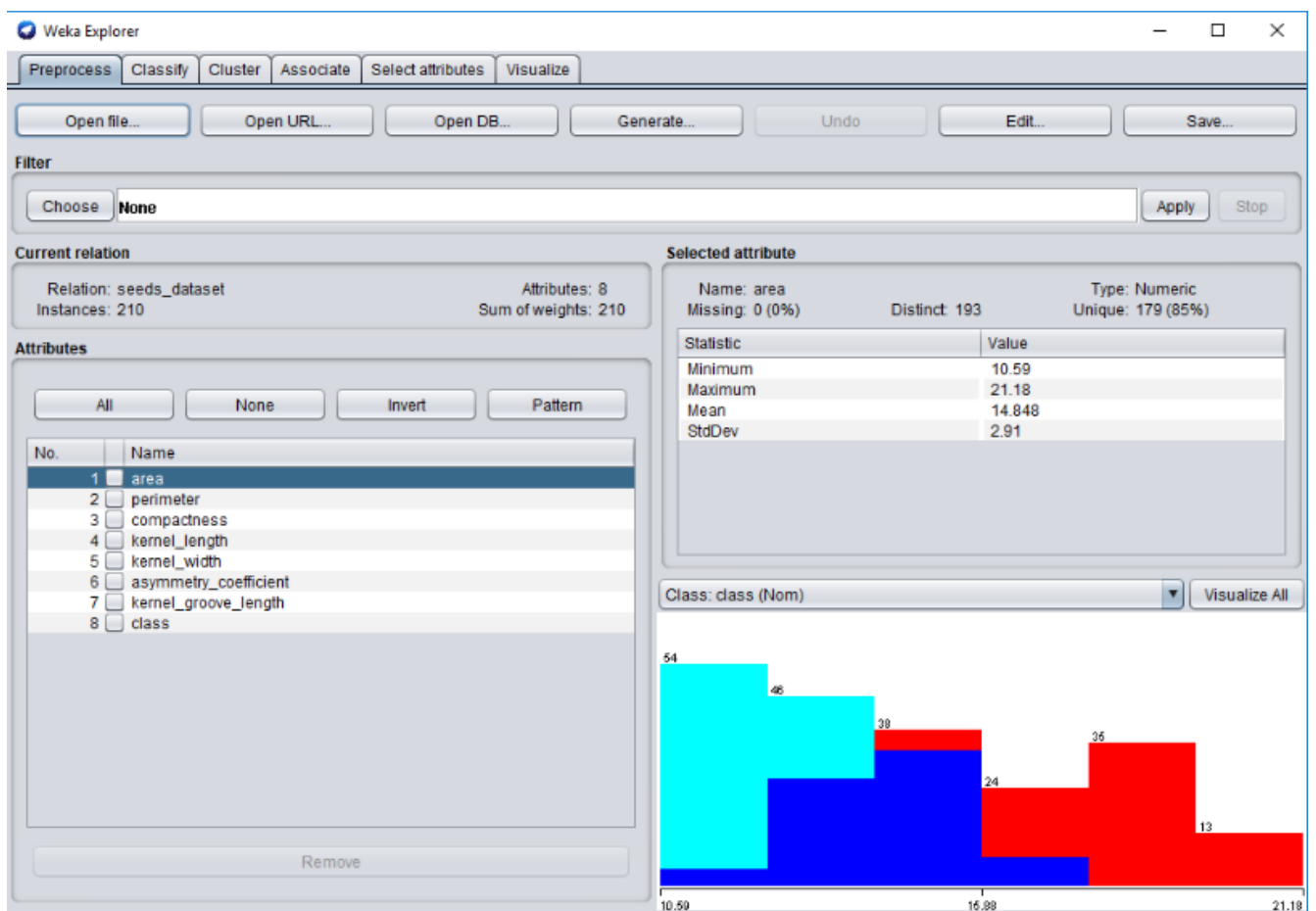


The data set did not contain attribute names and it had no missing values, but some data cells have shifted to the right as shown in the figure. All the values in the dataset were numeric.

After arranging the data set by making the data rows which have shifted to the right and putting the attribute names the dataset looks like below. In the class field, the values 1,2,3 were changed to Kama, Rosa and Canadian for proper visibility of association rules and clustering.

The dataset was then saved in .csv format and opened using Weka as shown below.

The csv file was then saved in .arff format.

The seeds_data.arff file was then loaded into weka and before doing the associate rule mining I did the clustering using k-means (SimpleKMeans), as after discretization clustering cannot be done.

As shown below the number of clusters were given as 3 and the Cluster mode was selected as Classes to clusters evaluation.



Then clicked OK and then clicked Start. Output of the k-means clustering can be seen below.

(Refer 'Kmeans-original.txt' for more details)

## Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

**Clusterer**

Choose | **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 3 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -

**Cluster mode**

- ○ Use training set
- ○ Supplied test set    Set...
- ○ Percentage split    % 66
- ● Classes to clusters evaluation
  - (Nom) class ▼
- ☑ Store clusters for visualization

Ignore attributes

Start | Stop

**Result list (right-click for options)**

17:46:15 - SimpleKMeans

**Clusterer output**

```
=== Run information ===

Scheme:       weka.clusterers.SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -
Relation:     seeds_dataset
Instances:    210
Attributes:   8
              area
              perimeter
              compactness
              kernel_length
              kernel_width
              asymmetry_coefficient
              kernel_groove_length
Ignored:
              class
Test mode:    Classes to clusters evaluation on training data


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


kMeans
======

Number of iterations: 5
Within cluster sum of squared errors: 22.024363075666038

Initial starting points (random):

Cluster 0: 18.59,16.05,0.9066,6.037,3.86,6.001,5.877
```

---

## Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

**Clusterer**

Choose | **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 3 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -

**Cluster mode**

- ○ Use training set
- ○ Supplied test set    Set...
- ○ Percentage split    % 66
- ● Classes to clusters evaluation
  - (Nom) class ▼
- ☑ Store clusters for visualization

Ignore attributes

Start | Stop

**Result list (right-click for options)**

17:46:15 - SimpleKMeans

**Clusterer output**

```
kMeans
======

Number of iterations: 5
Within cluster sum of squared errors: 22.024363075666038

Initial starting points (random):

Cluster 0: 18.59,16.05,0.9066,6.037,3.86,6.001,5.877
Cluster 1: 10.93,12.8,0.839,5.046,2.717,5.398,5.045
Cluster 2: 13.32,13.94,0.8613,5.541,3.073,7.035,5.44

Missing values globally replaced with mean/mode

Final cluster centroids:
                                    Cluster#
Attribute              Full Data        0        1        2
                        (210.0)     (64.0)   (77.0)   (69.0)
==============================================================
area                    14.8475    18.6102  11.8961  14.6512
perimeter               14.5593    16.2517  13.2577   14.442
compactness               0.871     0.8846   0.8498   0.8821
kernel_length            5.6285     6.1955   5.2306   5.5467
kernel_width             3.2586     3.7096    2.858   3.2873
asymmetry_coefficient    3.7002     3.5921   4.5995   2.7969
kernel_groove_length     5.4081     6.0567   5.0862   5.1656
```

**Status**

As seen from the figures the Incorrectly clustered instances are 23 and the percentage is 10.9524%.

Clustering can be visualized as shown below.

Then by ignoring attributes I tried to take the most accurate clustering with least incorrect instances.

Two results showed the minimum incorrect clustered instances as only 14 and the percentage as 6.6667%. One result was by ignoring the attributes kernel length and kernel width and the other result I got from ignoring area and kernel width as shown below.

Reasons for ignoring kernel width and kernel length in one result and then area and kernel width in another result was by looking at the point spread by visualize I identified according to their spread of points they should be ignored.

I tried for all other combinations as well to verify the minimum incorrectly clustered instances.

Result 1 (Ignoring kernel_length and kernel_width)

(Refer 'Kmeans-result1.txt' for more details)

## Weka Explorer

Preprocess | Classify | **Cluster** | Associate | Select attributes | Visualize

**Clusterer**

Choose | **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 3 -A"weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -

**Cluster mode**

- ○ Use training set
- ○ Supplied test set | Set...
- ○ Percentage split | % 66
- ● Classes to clusters evaluation
  - (Nom) class ▼
- ☑ Store clusters for visualization

Ignore attributes

Start | Stop

**Result list (right-click for options)**

17:46:15 - SimpleKMeans
17:56:08 - SimpleKMeans

**Clusterer output**

```
kMeans
======

Number of iterations: 6
Within cluster sum of squared errors: 17.031813652070092

Initial starting points (random):

Cluster 0: 18.59,16.05,0.9066,6.001,5.877
Cluster 1: 10.93,12.8,0.839,5.398,5.045
Cluster 2: 13.32,13.94,0.8613,7.035,5.44

Missing values globally replaced with mean/mode

Final cluster centroids:
                                    Cluster#
Attribute            Full Data          0          1          2
                       (210.0)     (69.0)     (71.0)     (70.0)
===============================================================
area                   14.8475    18.4078    11.8658    14.3624
perimeter              14.5593    16.1693    13.2565    14.2937
compactness              0.871     0.8836     0.8477     0.8822
asymmetry_coefficient   3.7002       3.63     4.7772      2.677
kernel_groove_length    5.4081     6.0358     5.1073     5.0944
```

---

## Weka Explorer

Preprocess | Classify | **Cluster** | Associate | Select attributes | Visualize

**Clusterer**

Choose | **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 3 -A"weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -

**Cluster mode**

- ○ Use training set
- ○ Supplied test set | Set...
- ○ Percentage split | % 66
- ● Classes to clusters evaluation
  - (Nom) class ▼
- ☑ Store clusters for visualization

Ignore attributes

Start | Stop

**Result list (right-click for options)**

17:46:15 - SimpleKMeans
17:56:08 - SimpleKMeans

**Clusterer output**

```
Time taken to build model (full training data) : 0.03 seconds

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

Clustered Instances

0      69 ( 33%)
1      71 ( 34%)
2      70 ( 33%)


Class attribute: class
Classes to Clusters:

  0  1  2  <-- assigned to cluster
  2  5 63 | Kama
 67  0  3 | Rosa
  0 66  4 | Canadian


Cluster 0 <-- Rosa
Cluster 1 <-- Canadian
Cluster 2 <-- Kama

Incorrectly clustered instances :      14.0      6.6667 %
```

## Result 2 (Ignoring area and kernel_width)

(Refer 'Kmeans-result2.txt' for more details)

## Accessing data directly from a database (JDBC)

I did also by using the JDBC connection method by creating a database in SQL server management studio and then by importing the csv file to that database as shown below.

Next JDBC driver should be downloaded and extracted. Then I edited the Database Utils.props by giving the JDBC connection database, username and password. Then enabled TCP/IP and enabled the 'sa' account by giving also a password and enabled SQL Authentication.

Then by using Open DB option in weka I connected to the SQL database and ran a SELECT * query to get all data from the data set.



Then clicked OK. This will direct to the Weka Explorer.

Weka Explorer screen view is given below.



All values including class comes in numeric format. Therefore, by using the weka filters by going to Choose->weka->filters->unsupervised->NumericToNominal I converted class attribute from numeric to nominal. Then I ran SimpleKMeans as same as done using the normal way which I did before I got the same inaccurate instances 14 and percentage as 10.9524%.

Then to make the dataset amenable for association rule mining all attributes should be converted to nominal. This was done by discretization as shown below.



I gave the number of bins as 3 and attributeIndices the attribute number (Parameter Setting)

Similarly, I gave number of bins as 3 after considering the amount of records in each attribute and got the output as below.



As labels are not meaningful, I had to edit using a text editor to make the values meaningful as shown below.



Then after opening saving it and opening it with weka the below window can be seen.

Now the data set is ready for the rule mining process.

# Rule mining process

For the rule mining process, I selected the Apriori algorithm.

## Reason for selecting Apriori algorithm

The Apriori Algorithm is considered an influential algorithm for mining frequent item sets for Boolean association rules. Apriori algorithm uses a "bottom up" approach, where frequent subsets are extended one item at a time (a step known as candidate generation) and then the groups of candidates are tested against the data.

The time required to run the Apriori algorithm after setting parameters for me was about 2 seconds. So Apriori Algorithm is fast.

Parameter Setting for the Apriori Algorithm was done as shown below.



I put the classindex as 8 because my 8th attribute is the class. I put lowerBoundMinSupport as 0.1 and upperBoundMinSupport, minMetric as 0.9 and firstly the number of rules were taken as 10 and the metricType as Confidence.

Later number of rules were taken as 1000 also to identify more rules.

Similarly, I changed the car parameter to True to get a more accurate and direct relationship with the class attribute which contains the 3 wheat varieties Kama, Rosa and Canadian.

Lift metricType was not taken because when car is True no rules were found using the lift metricType.

# Resulting rules

After running the Apriori algorithm the resulting rules for different parameter setting is given below.

1. Number of rules was taken as 10 and the metric type as confidence. Confidence level in all the 10 rules were 1 but since they are normally known they are not significant.13 cycles have been required to get the 10 rules where minMetric was 0.9.



(Refer 'Confidence,rules=10.txt' for the rules)

Confidence level in the rules varies from 0 to 1 and 1 is then taken as the best confidence.

2. Number of rules was taken as 1000 and the metric type as confidence.18 cycles have been required to get the 1000 rules where the minMetric was 0.9.

## Weka Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

**Associator**

Choose | Apriori -N 1000 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c 8

Start | Stop

**Associator output**

Result list (right-clic...

18:02:28 - Apriori
18:07:40 - Apriori
18:10:03 - Apriori

```
=== Run information ===

Scheme:       weka.associations.Apriori -N 1000 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c 8
Relation:     seeds_dataset-weka.filters.unsupervised.attribute.Discretize-B3-M-1.0-R1-precision6-weka.filters.unsupe
Instances:    210
Attributes:   8
              area
              perimeter
              compactness
              kernel_length
              kernel_width
              asymmetry_coefficient
              kernel_groove_length
              class
=== Associator model (full training set) ===
```

---

## Weka Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

**Associator**

Choose | Apriori -N 1000 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c 8

Start | Stop

**Associator output**

Result list (right-clic...

18:02:28 - Apriori
18:07:40 - Apriori
18:10:03 - Apriori

```
Apriori
=======

Minimum support: 0.1 (21 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 18

Generated sets of large itemsets:

Size of set of large itemsets L(1): 23

Size of set of large itemsets L(2): 128

Size of set of large itemsets L(3): 215

Size of set of large itemsets L(4): 191

Size of set of large itemsets L(5): 105

Size of set of large itemsets L(6): 39

Size of set of large itemsets L(7): 9

Size of set of large itemsets L(8): 1

Best rules found:

 1. perimeter=below_14.023333 90 ==> area=below_14.12 90    <conf:(1)> lift:(2.1) lev:(0.22) [47] conv:(47.14)
 2. perimeter=below_14.023333 kernel_length=below_5.491 87 ==> area=below_14.12 87    <conf:(1)> lift:(2.1) lev:(0.2
```

```
Weka Explorer                                                               —  □  ×

 Preprocess  Classify  Cluster  Associate  Select attributes  Visualize
Associator

  Choose   Apriori -N 1000 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c 8

                     Associator output
  Start    Stop      974. area=below_14.12 perimeter=below_14.023333 compactness=below_0.844833 kernel_length=below_5.491 kernel_width=bel
Result list (right-clic...  975. compactness=below_0.844833 kernel_length=below_5.491 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-
                     976. perimeter=below_14.023333 compactness=below_0.844833 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-
  18:02:28 - Apriori 977. perimeter=below_14.023333 compactness=below_0.844833 kernel_length=below_5.491 asymmetry_coefficient=3.328733-5.
  18:07:40 - Apriori 978. perimeter=below_14.023333 compactness=below_0.844833 kernel_length=below_5.491 kernel_width=below_3.097667 asymm
  18:10:03 - Apriori 979. area=below_14.12 compactness=below_0.844833 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-5.892367
                     980. area=below_14.12 compactness=below_0.844833 kernel_length=below_5.491 asymmetry_coefficient=3.328733-5.892367 cl
                     981. area=below_14.12 compactness=below_0.844833 kernel_length=below_5.491 kernel_width=below_3.097667 asymmetry_coef
                     982. area=below_14.12 perimeter=below_14.023333 compactness=below_0.844833 asymmetry_coefficient=3.328733-5.892367 cl
                     983. area=below_14.12 perimeter=below_14.023333 compactness=below_0.844833 kernel_width=below_3.097667 asymmetry_coef
                     984. area=below_14.12 perimeter=below_14.023333 compactness=below_0.844833 kernel_length=below_5.491 asymmetry_coeffi
                     985. compactness=below_0.844833 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-5.892367 class=Canadian 23
                     986. compactness=below_0.844833 kernel_length=below_5.491 asymmetry_coefficient=3.328733-5.892367 class=Canadian 23 =
                     987. compactness=below_0.844833 kernel_length=below_5.491 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-
                     988. perimeter=below_14.023333 compactness=below_0.844833 asymmetry_coefficient=3.328733-5.892367 class=Canadian 23 =
                     989. perimeter=below_14.023333 compactness=below_0.844833 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-
                     990. perimeter=below_14.023333 compactness=below_0.844833 kernel_length=below_5.491 asymmetry_coefficient=3.328733-5.
                     991. area=below_14.12 compactness=below_0.844833 asymmetry_coefficient=3.328733-5.892367 class=Canadian 23 ==> perime
                     992. area=below_14.12 compactness=below_0.844833 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-5.892367
                     993. area=below_14.12 compactness=below_0.844833 kernel_length=below_5.491 asymmetry_coefficient=3.328733-5.892367 23
                     994. area=below_14.12 perimeter=below_14.023333 compactness=below_0.844833 asymmetry_coefficient=3.328733-5.892367 23
                     995. compactness=below_0.844833 asymmetry_coefficient=3.328733-5.892367 class=Canadian 23 ==> area=below_14.12 perime
                     996. compactness=below_0.844833 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-5.892367 23 ==> area=below
                     997. compactness=below_0.844833 kernel_length=below_5.491 asymmetry_coefficient=3.328733-5.892367 23 ==> area=below_1
                     998. perimeter=below_14.023333 compactness=below_0.844833 asymmetry_coefficient=3.328733-5.892367 23 ==> area=below_1
                     999. area=below_14.12 compactness=below_0.844833 asymmetry_coefficient=3.328733-5.892367 23 ==> perimeter=below_14.02
                     1000. compactness=below_0.844833 asymmetry_coefficient=3.328733-5.892367 23 ==> area=below_14.12 perimeter=below_14.0
```

(Refer 'Confidence,rules=1000.txt' for the rules)

Some interesting rule descriptions found are

i)      If perimeter is above 15.636667 then (==>) the wheat variety would be Rosa
        (class=Rosa). This was shown with a confidence of 1.

ii)     If perimeter is above 15.636667 and kernel width is above 3.565333 48 then (==>)
        the wheat variety would be Rosa (class=Rosa). This was shown with a confidence of
        1.

iii)    If perimeter is between 14.023333 and 15.636667and asymmetry coefficient is
        below3.328733 and kernel groove length is below_5.196   then (==>) kernel width is
        between 3.097667 and 3.565333 and the wheat variety would be Kama
        (class=Kama). This was shown with a confidence of 1.

iv)     If area is below 14.12 and perimeter is below14.023333 and compactness is below
        0.844833 and kernel length is below 5.491 and kernel width is below 3.097667 and
        asymmetry coefficient is between 3.328733 and 5.892367 then (==>) wheat variety
        is Canadian (class=Canadian). This was shown with a confidence of 1.

(All the attributes are the geometrical properties of kernels belonging to the three different
varieties of wheat namely Kama, Rosa and Canadian)

All the above were done taking the car as False. More direct relationships can be taken by taking car as True so that the right side will have the class attribute.

3. When the car parameter is taken as true and the number of rules is given as 1000 and the metricType as confidence we can see 155 rules as seen below. The number of cycles is 18 and the minMetric is 0.9.

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

**Associator**

Choose | Apriori -N 1000 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -A -c 8

**Associator output**

```
129. perimeter=14.023333-15.636667 kernel_width=3.097667-3.565333 asymmetry_coefficient=below_3.328733 41 ==> class
130. area=14.12_17.65 kernel_width=3.097667-3.565333 asymmetry_coefficient=below_3.328733 32 ==> class=Kama 30      c
131. area=14.12_17.65 perimeter=14.023333-15.636667 kernel_width=3.097667-3.565333 asymmetry_coefficient=below_3.32
132. kernel_length=5.491-6.083 kernel_width=3.097667-3.565333 asymmetry_coefficient=below_3.328733 31 ==> class=Kam
133. kernel_width=3.097667-3.565333 asymmetry_coefficient=below_3.328733 kernel_groove_length=below_5.196 31 ==> cl
134. perimeter=14.023333-15.636667 kernel_length=5.491-6.083 kernel_width=3.097667-3.565333 asymmetry_coefficient=b
135. perimeter=14.023333-15.636667 asymmetry_coefficient=below_3.328733 43 ==> class=Kama 40      conf:(0.93)
136. compactness=above_0.881567 kernel_width=3.097667-3.565333 asymmetry_coefficient=below_3.328733 27 ==> class=Ka
137. area=14.12_17.65 kernel_length=5.491-6.083 kernel_width=3.097667-3.565333 asymmetry_coefficient=below_3.328733
138. area=14.12_17.65 perimeter=14.023333-15.636667 kernel_length=5.491-6.083 kernel_width=3.097667-3.565333 asymme
139. kernel_length=below_5.491 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-5.892367 53 ==> class=Can
140. area=below_14.12 kernel_length=below_5.491 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-5.892367
141. perimeter=below_14.023333 kernel_length=below_5.491 kernel_width=below_3.097667 asymmetry_coefficient=3.328733
142. area=below_14.12 perimeter=below_14.023333 kernel_length=below_5.491 kernel_width=below_3.097667 asymmetry_coe
143. area=below_14.12 kernel_groove_length=5.196-5.873 23 ==> class=Canadian 21      conf:(0.91)
144. compactness=above_0.881567 asymmetry_coefficient=below_3.328733 kernel_groove_length=below_5.196 23 ==> class=
145. area=14.12_17.65 asymmetry_coefficient=below_3.328733 34 ==> class=Kama 31      conf:(0.91)
146. area=14.12_17.65 perimeter=14.023333-15.636667 asymmetry_coefficient=below_3.328733 34 ==> class=Kama 31      co
147. perimeter=14.023333-15.636667 kernel_length=5.491-6.083 asymmetry_coefficient=below_3.328733 33 ==> class=Kama
148. kernel_width=below_3.097667 asymmetry_coefficient=3.328733-5.892367 54 ==> class=Canadian 49      conf:(0.91)
149. area=below_14.12 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-5.892367 54 ==> class=Canadian 49
150. perimeter=below_14.023333 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-5.892367 54 ==> class=Can
151. area=below_14.12 perimeter=below_14.023333 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-5.892367
152. kernel_length=below_5.491 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-5.892367 kernel_groove_le
153. area=below_14.12 kernel_length=below_5.491 kernel_width=below_3.097667 asymmetry_coefficient=3.328733-5.892367
154. perimeter=below_14.023333 kernel_length=below_5.491 kernel_width=below_3.097667 asymmetry_coefficient=3.328733
155. area=below_14.12 perimeter=below_14.023333 kernel_length=below_5.491 kernel_width=below_3.097667 asymmetry_coe
```

Start | Stop

Result list (right-clic...

18:02:28 - Apriori
18:07:40 - Apriori
18:10:03 - Apriori
18:36:26 - Apriori

All these 155 rules are very interesing rules and this 155 rules will be the selection which would be shown to the client. From the 155 rules the top or the first 104 rules are the best rules with confidence 1. So the best selection to show for the client is this 104 rules. The 155 rules is attached with this report which is in a text file. (Refer 'Confidence,Car=true.txt' for the rules)

Some one to one or simple relationship rule descriptions from the 155 rules are given below.
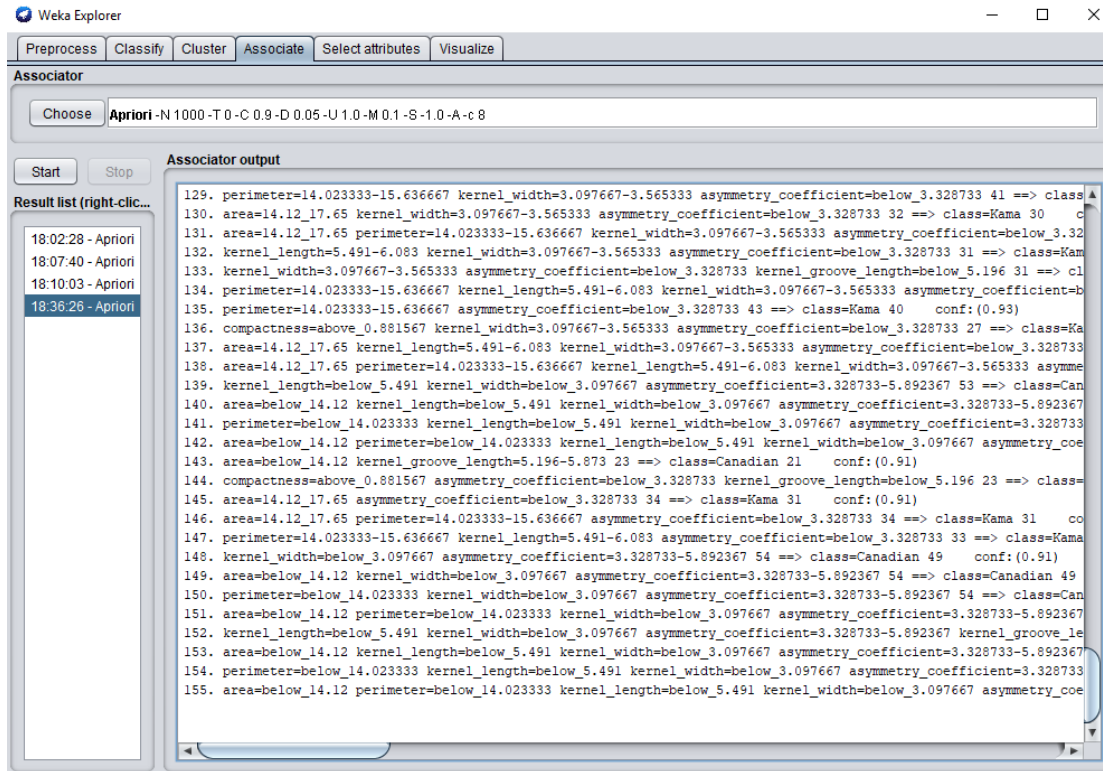
i.   If perimeter is above 15.636667  then (==>) the wheat variety would be Rosa (class=Rosa). This was shown with a confidence of 1.

ii.  If area is above 17.65 then(==>)  the wheat variety would be Rosa (class=Rosa). This was shown with a confidence of 1.

iii. If compactness is below 0.844833  then (==>) the wheat variety would be Canadian (class=Canadian). This was shown with a confidence of 0.97.

iv.  If area is 14.12_17.65 and asymmetry coefficient is below 3.328733 then (==>) the wheat variety would be Kama (class=Kama). This was shown with a confidence of 0.91.

v.   If area is below 14.12 and kernel groove length is between 5.196 and 5.873 then (==>) the wheat variety would be Canadian (class=Canadian). This was shown with a confidence of 0.91.

(All the attributes are the geometrical properties of kernels belonging to the three different varieties of wheat namely Kama, Rosa and Canadian)

# Recommendations

Mainly to identify the 3 wheat varieties Kama, Rosa and Canadian I can recommend 3 recommendations considering each wheat variety to the client.

1) If perimeter of the kernel is greater than 15.64 and area of the kernel is greater than 17.65 then I would recommend considering the wheat variety as Rosa.
2) If kernel width is below 3.10 and kernel groove length is between 5.20 and 5.87, I would recommend considering the wheat variety as Canadian.
3) If perimeter is between 14.02 and 15.64 and asymmetry coefficient is below 3.33 and kernel groove length is below 5.20, I would be recommend considering the wheat variety as Canadian or perimeter is between 14.02 and 15.64 and kernel width is between 3.10 and 3.57 and kernel groove length is below 5.20, I would recommend considering the wheat variety as Canadian.

(All values are taken after rounding to the second decimal place)

Above 3 recommendations were taken after considering the rules with confidence equal to 1.

From above recommendations the client can simple identify the 3 wheat varieties Kama, Rosa and Canadian separately.