

Shri Ramdeobaba College of Engineering and Management, Nagpur
Department of Computer Science and Engineering
Session: 2021-2022

Name-Harsh Agrawal CSE SEM-7 Section-B Shift-2

Data Visualization and Analytics Lab

VII Semester

PRACTICAL NO. 1

Aim: Introduction to weka and data preprocessing on the given data set in Weka

Theory Questions:

1. What options are available on the main panel?

- Explorer
- Experimenter
- KnowledgeFlow
- Workbench
- Simple CLI

2. What is the purpose of the the following in Weka:

1. The Explorer

→ The WEKA Explorer windows show different tabs starting with preprocessing. Initially, the preprocess tab is active, as first the data set is preprocessed before applying algorithms to it and exploring the dataset. It is a platform to work with data and apply the transformation. We can effectively work with different data mining and machine learning algorithms.

2. The Knowledge Flow interface

→ Knowledge flow shows a graphical representation of WEKA algorithms. The user can select the components and create a workflow to analyze the datasets. The data can be handled batch-wise or incrementally. Parallel workflows can be designed and each will run in a separate thread.

3. The Experimenter

→ The WEKA experimenter button allows the users to create, run, and modify different schemes in one experiment on a dataset. The experimenter has 2 types of configuration: Simple and Advanced. Both configurations allow users to run experiments locally and on remote computers.

4. The command-line interface

→ Simple CLI is Weka Shell with command line and output. With “help”, the overview of all the commands can be seen. Simple CLI offers access to all classes such as classifiers, clusters, filters, etc.

3. Describe the arff file format.

→ WEKA works on the ARFF file for data analysis. ARFF stands for Attribute Relation File Format. It has 3 sections: relation, attributes, and data. Every section starts with “@”.

ARFF files take Nominal, Numeric, String, Date, and Relational data attributes. Some of the well-known machine learning datasets are present in WEKA as ARFF.

Ex:- @relation new

@attribute outlook {sunny, overcast, rainy}

@attribute temp {hot, humid}

@attribute age real

@attribute play {yes, no}

@data

sunny, hot, 12, yes

rainy, hot, 45, no

4. What is the purpose of the following in the Explorer Panel?

1. The Preprocess panel

→ This allows us to choose the data file.

1. What are the main sections of the Preprocess panel?

2. What are the primary sources of data in Weka?

→ Open File – enables the user to select the file from the local machine

→ Open URL – enables the user to select the data file from different locations

→ Open Database – enables users to retrieve a data file from a database source

2. The Classify panel

→ The classifier panel allows you to configure and execute any of the weka classifiers on the current dataset. You can choose to perform a cross validation or

test on a separate dataset. Classification errors can be visualized in a pop-up data visualization tool. If the classifier produces a decision tree it can be displayed graphically in a pop-up tree visualizer.

3. The Cluster panel

→ From the cluster panel you can configure and execute any of the weka clusterers on the current dataset. Clusters can be visualized in a pop-up data visualization tool.

4. The Associate panel

→ This allows us to apply association rules, which identify the association within the data. Algorithms can be association rule mining and apriori.

5. The Select Attributes panel

→ This panel allows you to configure and apply any combination of weka attribute evaluator and search method to select the most pertinent attributes in the dataset. If an attribute selection scheme transforms the data then the transformed data can be visualized in a pop-up data visualization tool.

6. The Visualize panel.

→ This panel displays a scatter plot matrix for the current dataset. The size of the individual cells and the size of the points they display can be adjusted using the slider controls at the bottom of the panel. It lets us look at the dataset and select different attributes, preferably numeric ones from x and y axes.

Dataset Link [Weather.nominal.arff]:

<https://storm.cis.fordham.edu/~gweiss/data-mining/weka-data/weather.nominal.arff>

Dataset Link [sick.arff]:

https://datahub.io/machine-learning/sick#resource-sick_arff

Questions:

1. Press the Explorer button on the main panel and load the weather dataset and answer the following questions

1. How many instances are there in the dataset?

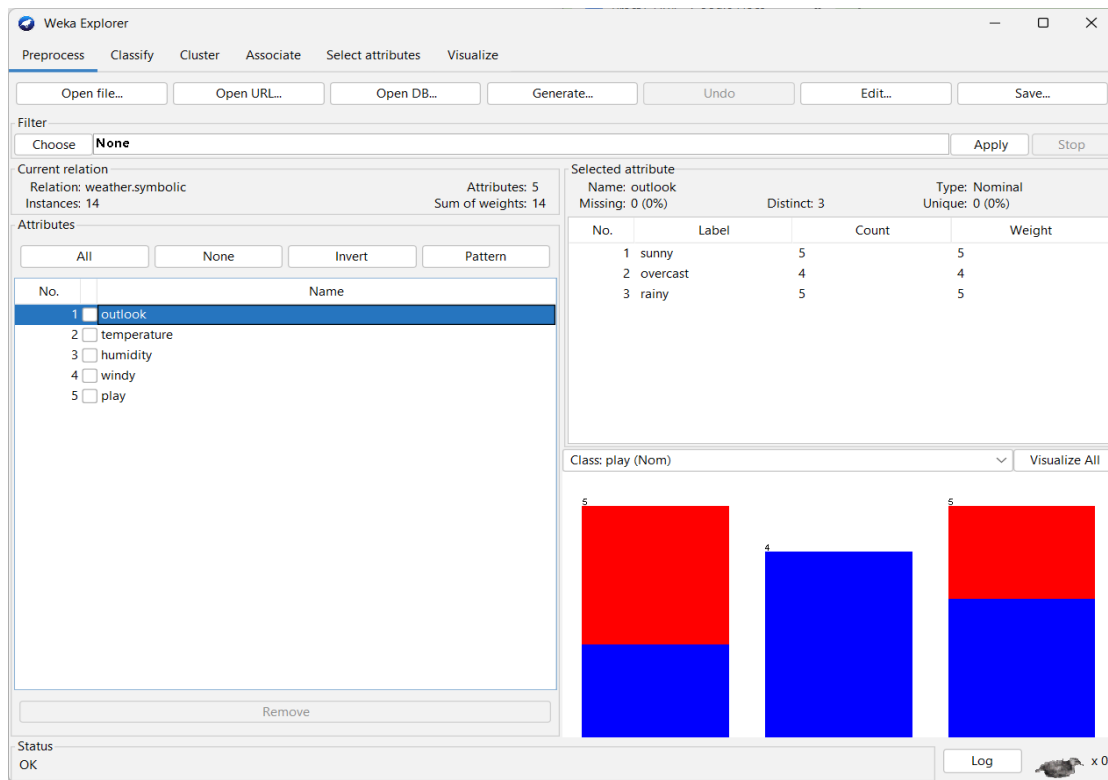
=> 14

Viewer					
Relation: weather.symbolic					
No.	1: outlook Nominal	2: temperature Nominal	3: humidity Nominal	4: windy Nominal	5: play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no

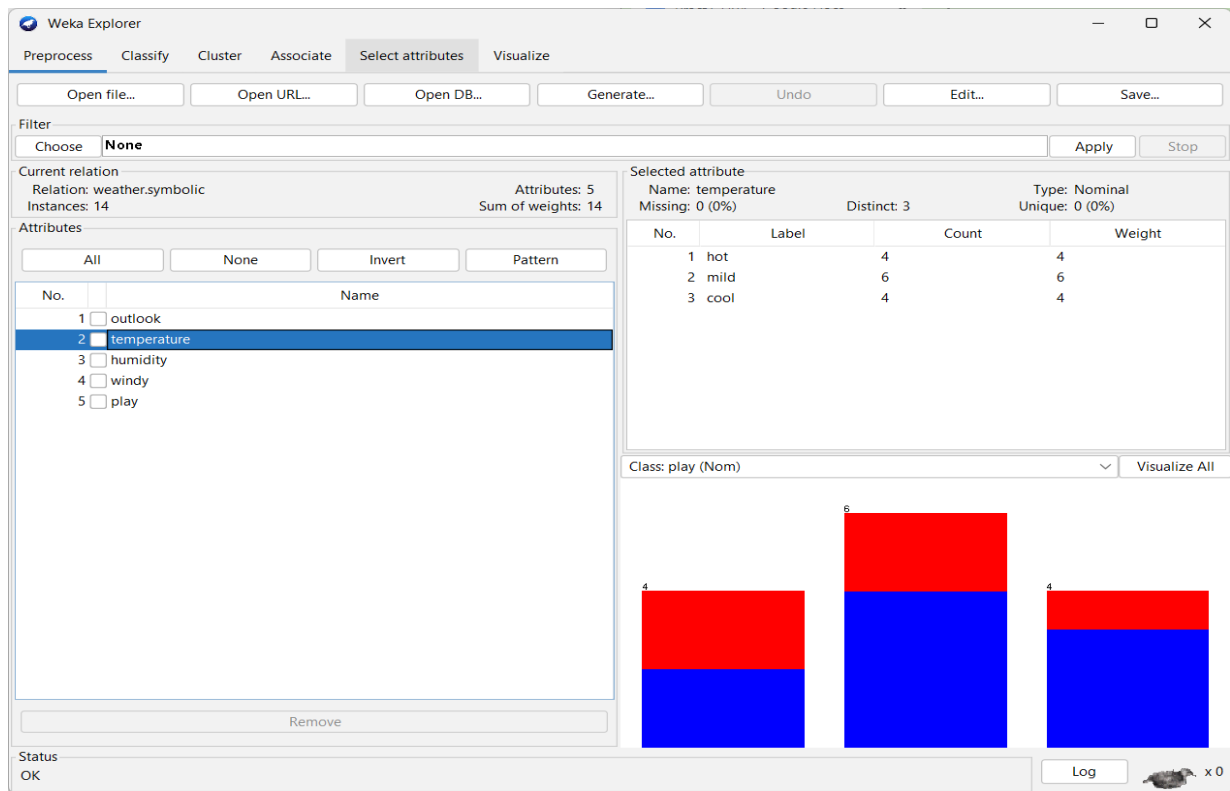
2. State the names of the attributes along with their types and values.

→ outlook {sunny, overcast, rainy}, type nominal
 temperature {hot, mild, cool}, type nominal
 humidity {high, normal}, type nominal
 windy {TRUE, FALSE}, type nominal
 play {yes, no}, type nominal

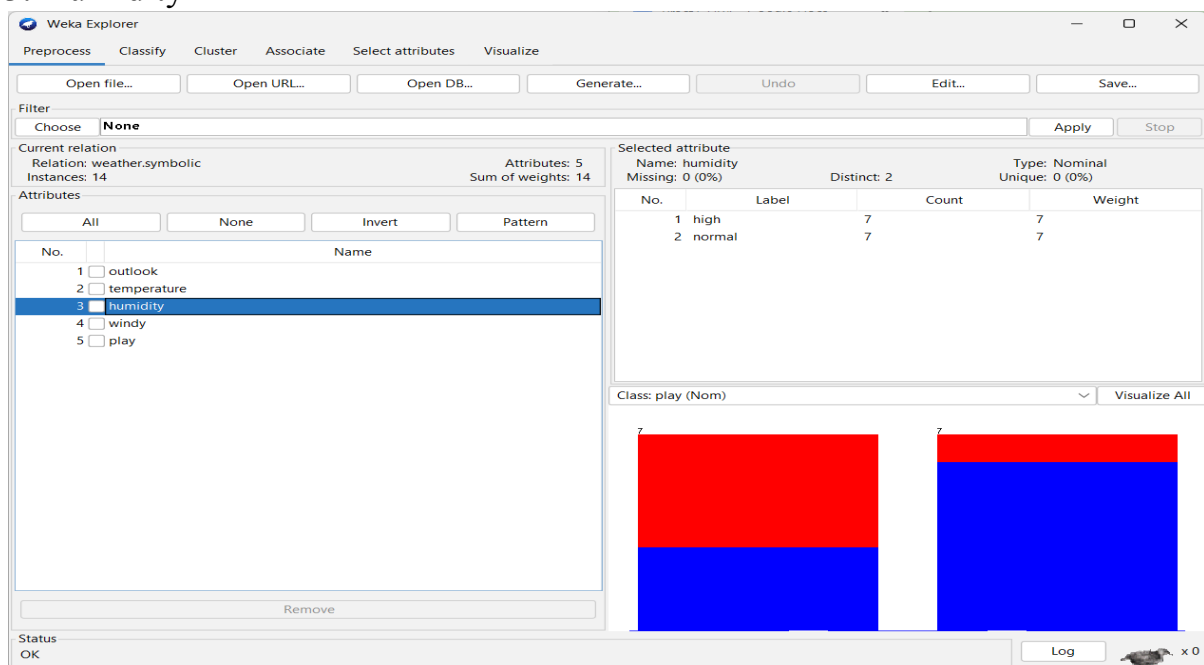
1. Outlook



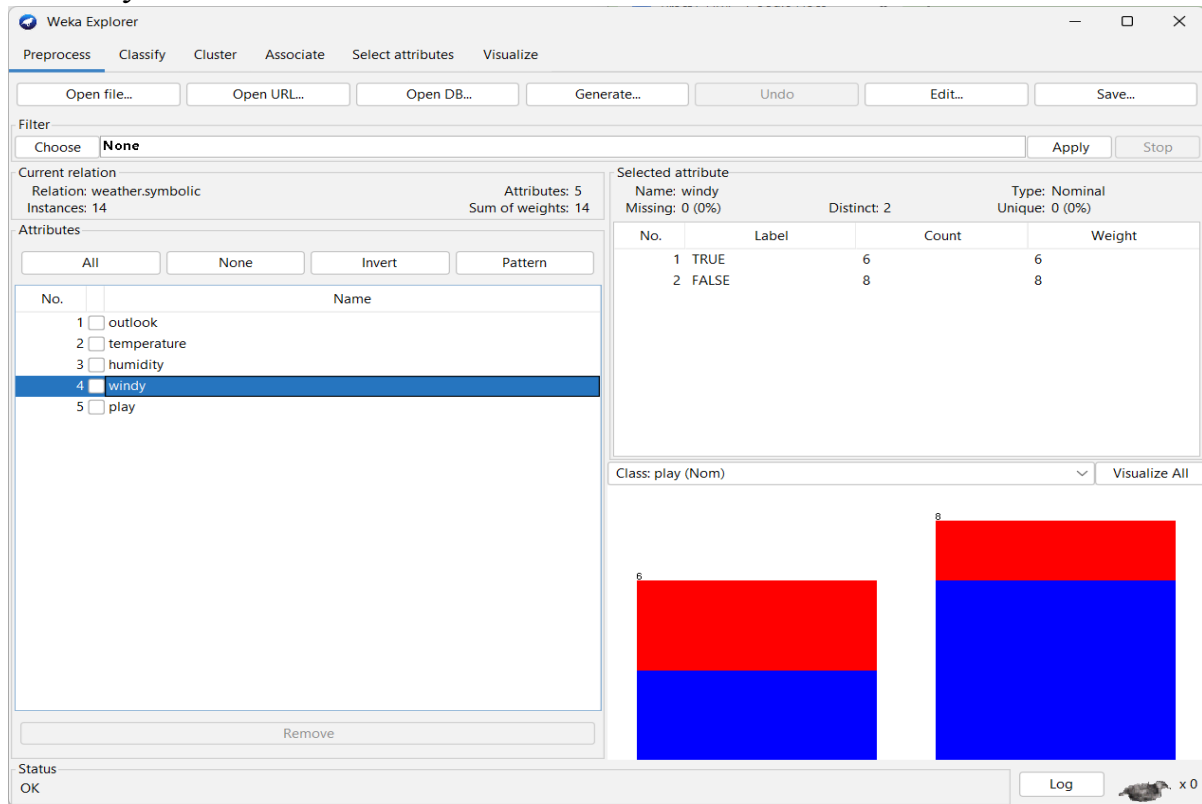
2. Temperature



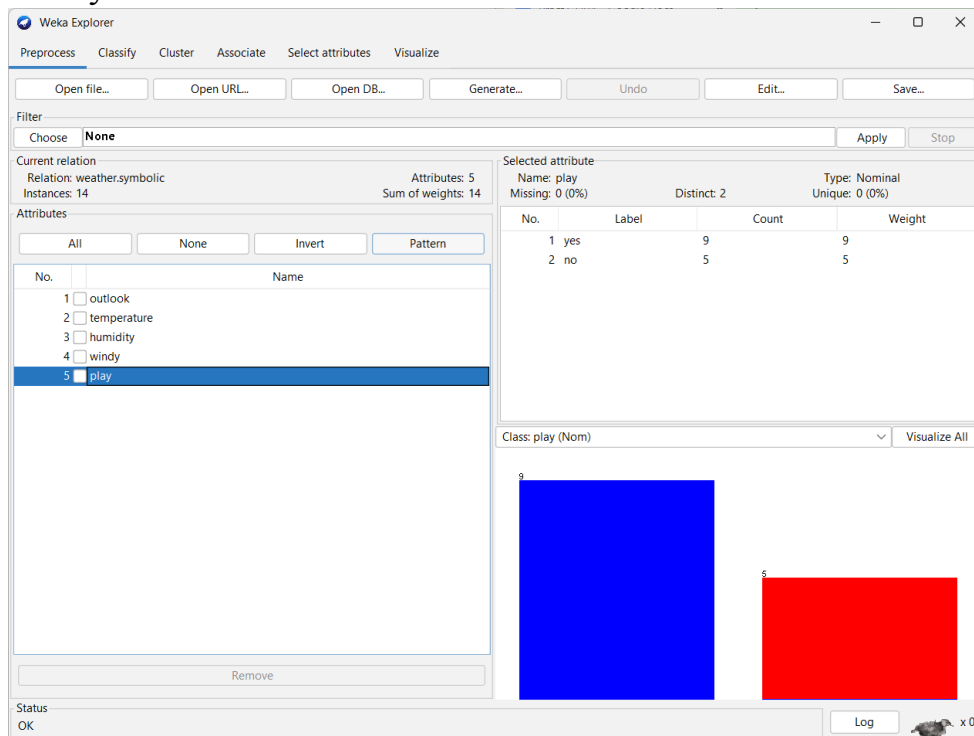
3. Humidity



4. Windy



5. Play



3. What is the class attribute?

→ play

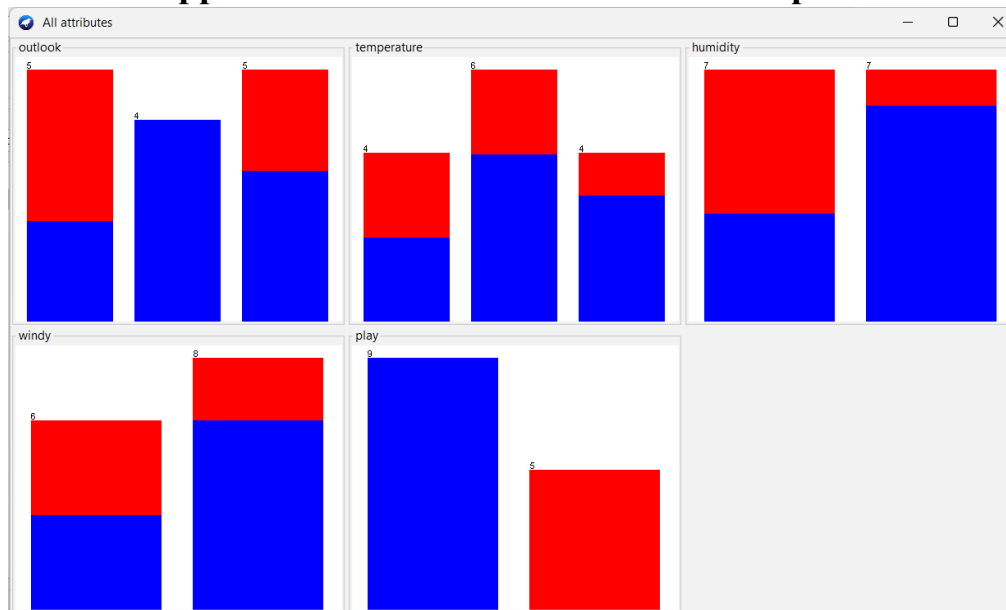
4. How will you determine how many instances of each class are present in the data

→ Click on edit and see the rows

We can hover the histogram to determine the number of instances of each class present in the data.

Viewer					
Relation: weather.symbolic					
No.	1: outlook Nominal	2: temperature Nominal	3: humidity Nominal	4: windy Nominal	5: play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no

5. What happens when the Visualize All button is pressed?



6. How will you view the instances in the dataset? How will you save the changes?

To view instances in dataset we can select button and view all the instances and then select various options given below to perform respective operations. OK to save it.

Viewer

Relation: weather.symbolic

No.	1: outlook Nominal	2: temperature Nominal	3: humidity Nominal	4: windy Nominal	5: play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no
15	sunny	cool	high	TRUE	yes

hot
mild
cool

Add instance Undo OK Cancel

7. Now, extend the dataset to include 50 instances in total.

Viewer

Relation: weather.symbolic

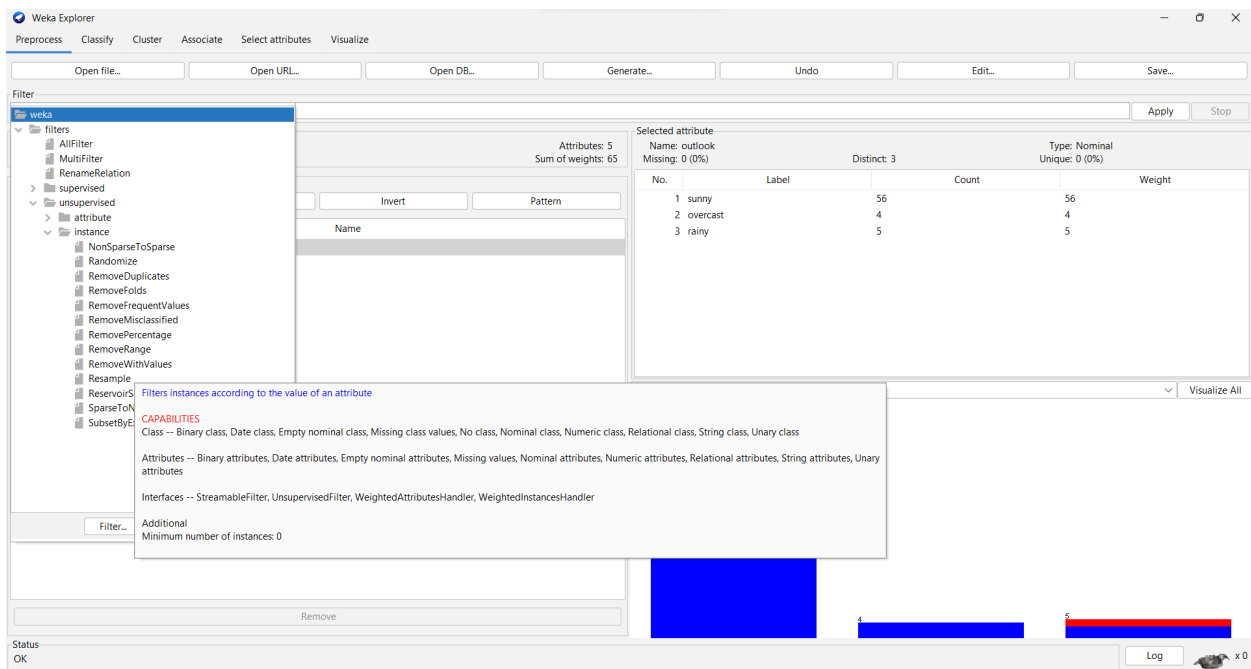
No.	1: outlook Nominal	2: temperature Nominal	3: humidity Nominal	4: windy Nominal	5: play Nominal
27	sunny	hot	high	TRUE	yes
28	sunny	hot	high	TRUE	yes
29	sunny	hot	high	TRUE	yes
30	sunny	hot	high	TRUE	yes
31	sunny	hot	high	TRUE	yes
32	sunny	hot	high	TRUE	yes
33	sunny	hot	high	TRUE	yes
34	sunny	hot	high	TRUE	yes
35	sunny	hot	high	TRUE	yes
36	sunny	hot	high	TRUE	yes
37	sunny	hot	high	TRUE	yes
38	sunny	hot	high	TRUE	yes
39	sunny	hot	high	TRUE	yes
40	sunny	hot	high	TRUE	yes
41	sunny	hot	high	TRUE	yes
42	sunny	hot	high	TRUE	yes
43	sunny	hot	high	TRUE	yes
44	sunny	hot	high	TRUE	yes
45	sunny	hot	high	TRUE	yes
46	sunny	hot	high	TRUE	yes
47	sunny	hot	high	TRUE	yes
48	sunny	hot	high	TRUE	yes
49	sunny	hot	high	TRUE	yes
50	sunny	hot	high	TRUE	yes

Add instance Undo OK Cancel

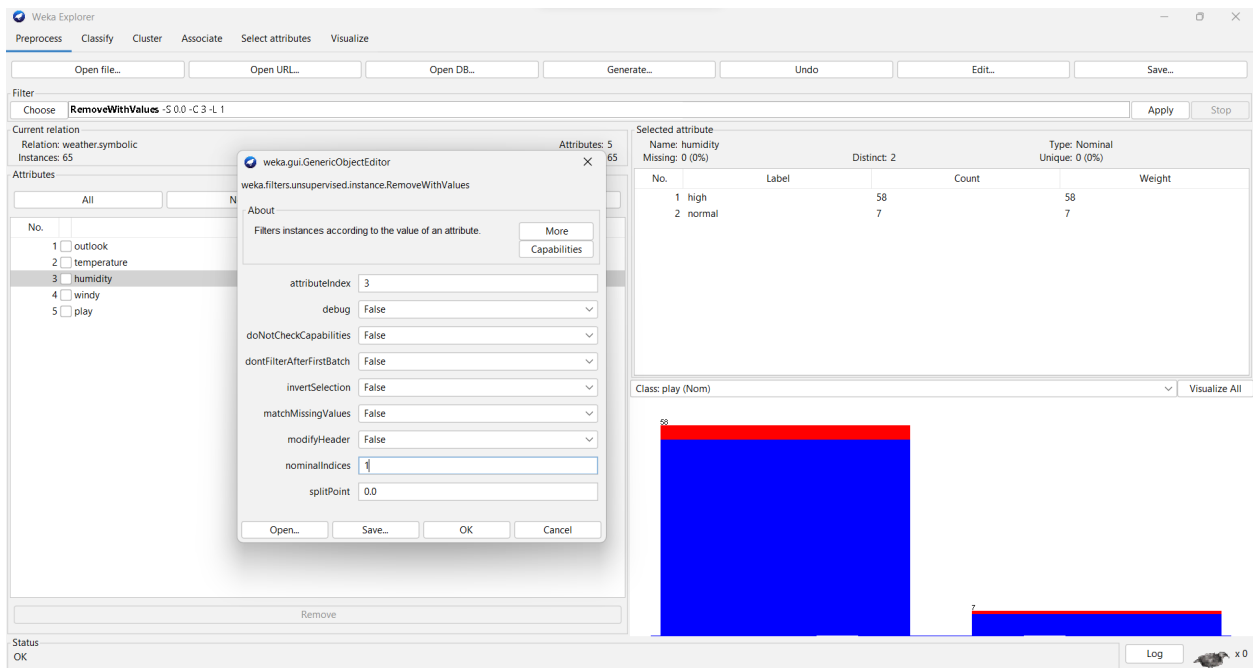
2. Do as directed to apply Filter

1. Use the unsupervised filter RemoveWithValues to remove all instances where the attribute 'humidity' has the value 'high'? Undo the effect of the filter.

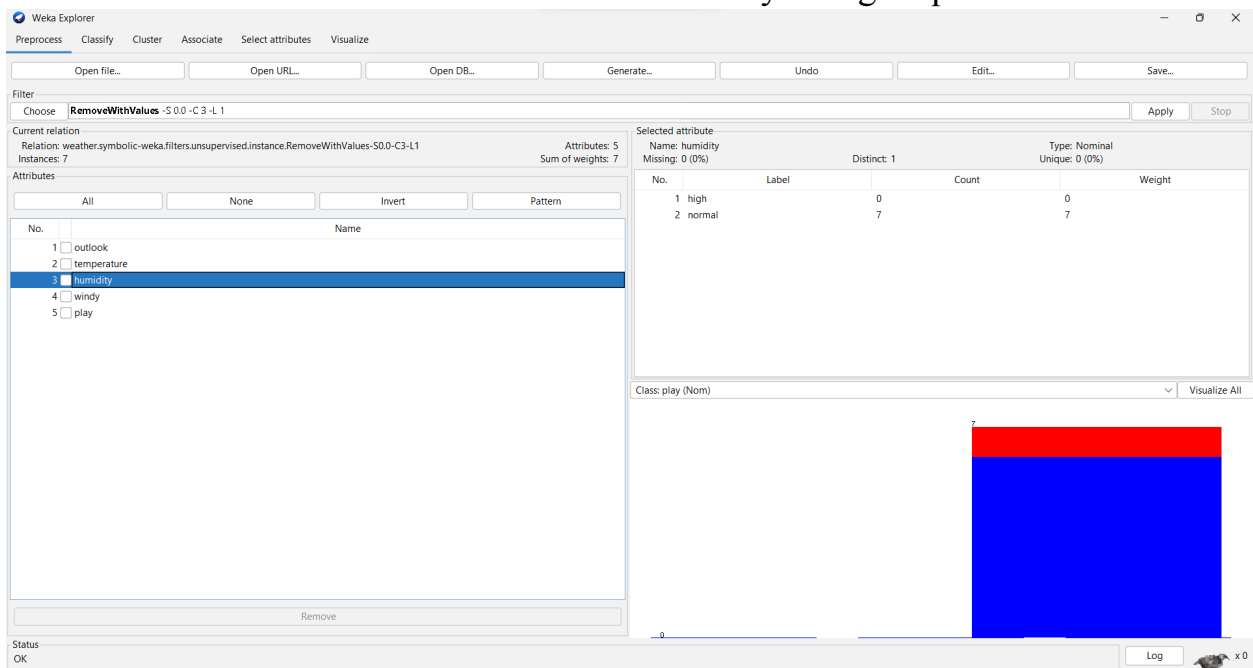
→ From the weka filter options Remove WithValues filter from the unsupervised and then the instance category. Here by applying the filter we can mark the corrupt values as missing in the dataset, remove instances with missing values in the dataset etc.



Attribute humidity is selected.



All the instance values where the humidity value is high has been removed from the dataset. We can see the count value for humidity as high equal to zero.



→ UNDO button has been pressed to undo all the applied effects.

2. Remove the 'FALSE' instances of windy attribute and undo the effect.

Weka Explorer interface showing the 'RemoveWithValues' filter applied to the 'windy' attribute. The filter dialog is open, showing 'attributeIndex: 4' and 'splitPoint: 0.0'. The 'Selected attribute' table shows 'TRUE' with count 57 and 'FALSE' with count 8. The bar chart shows a large blue bar for 'TRUE' and a small red bar for 'FALSE'.

No.	Label	Count	Weight
1	TRUE	57	57
2	FALSE	8	8

→ Similar to as done above, all the false instances of the windy attribute has been removed.

3. Remove the attribute outlook and undo the effect.

Weka Explorer interface showing the 'RemoveDuplicates' filter applied to the 'outlook' attribute. The filter dialog is open, showing 'attributeIndex: 1' and 'splitPoint: 0.0'. The 'Selected attribute' table shows 'sunny' with count 6, 'overcast' with count 4, and 'rainy' with count 5. The bar chart shows three bars: 'sunny' (blue), 'overcast' (red), and 'rainy' (blue).

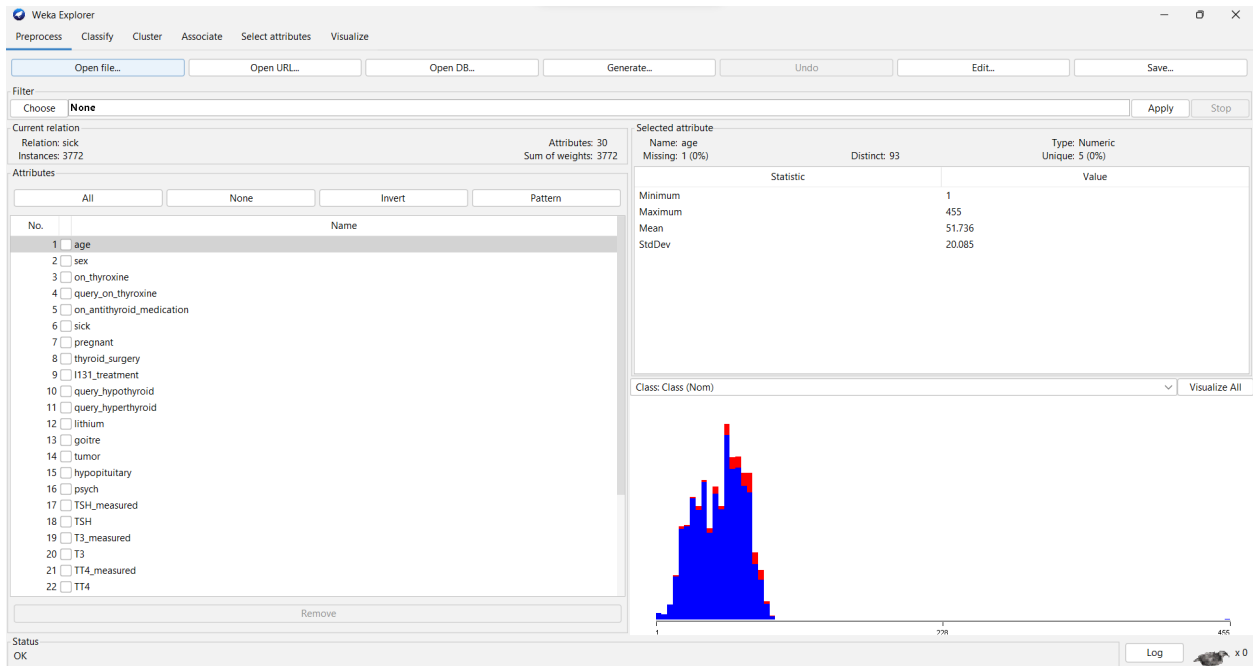
No.	Label	Count	Weight
1	sunny	6	6
2	overcast	4	4
3	rainy	5	5

4. Experiment with different filters and report their effects.

3. Application of Discretization Filters [use sick.arff dataset]

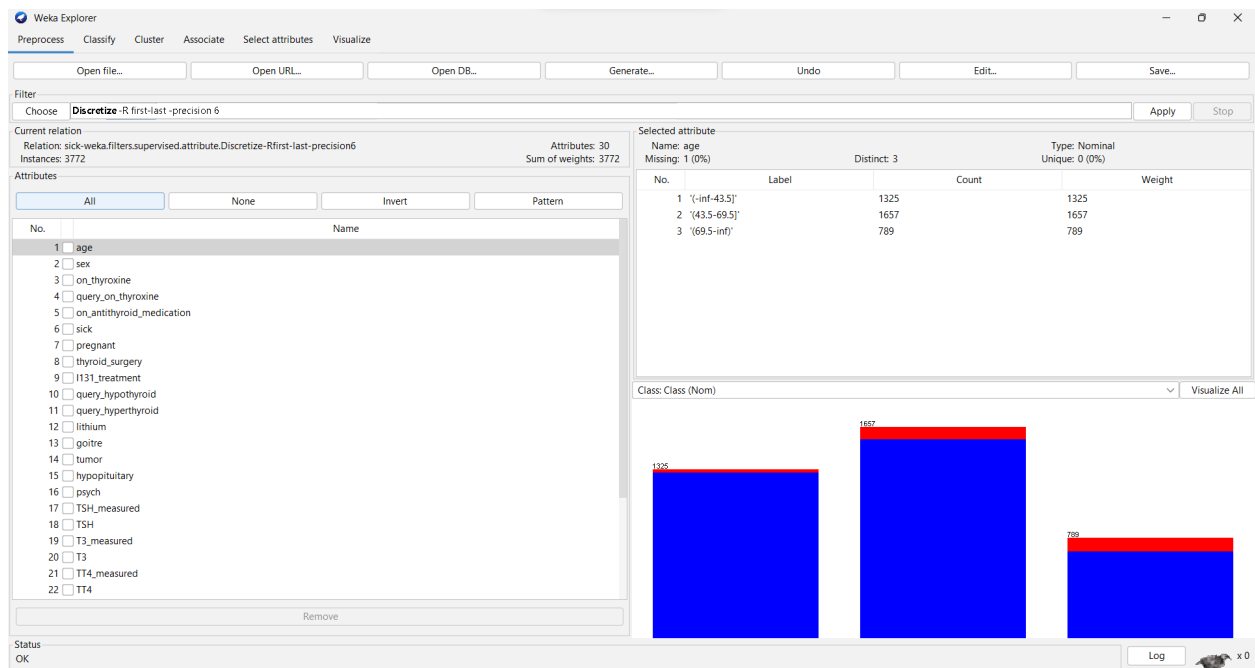
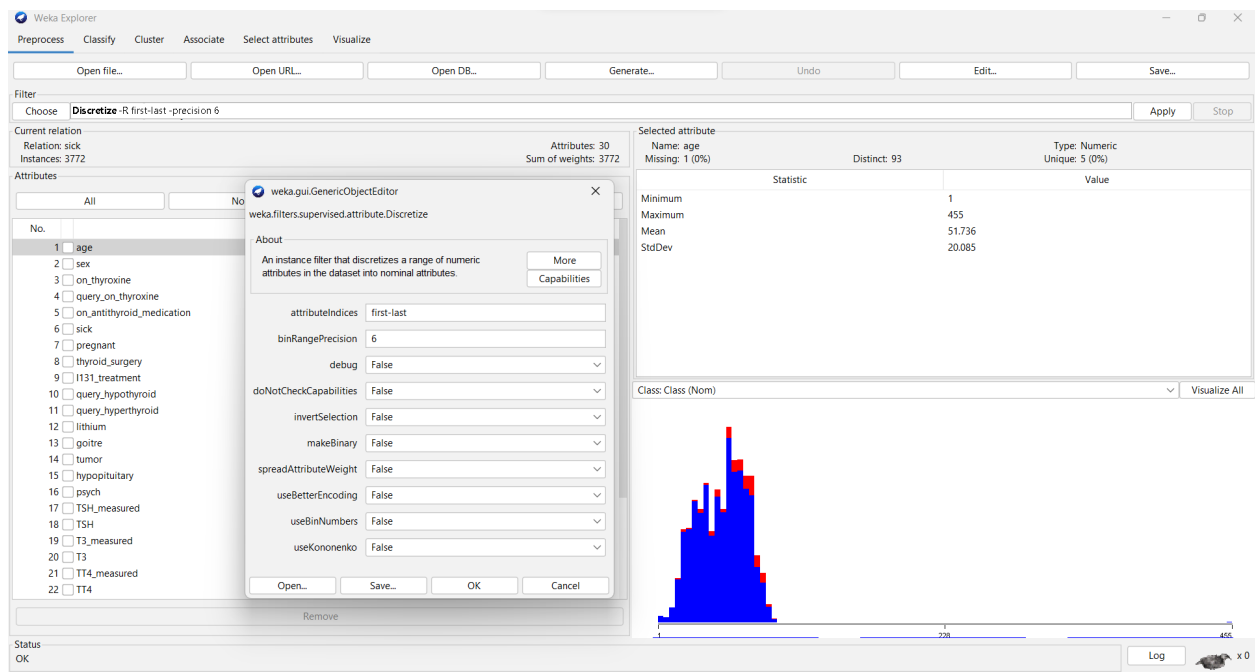
1. Load the 'sick.arff' dataset.

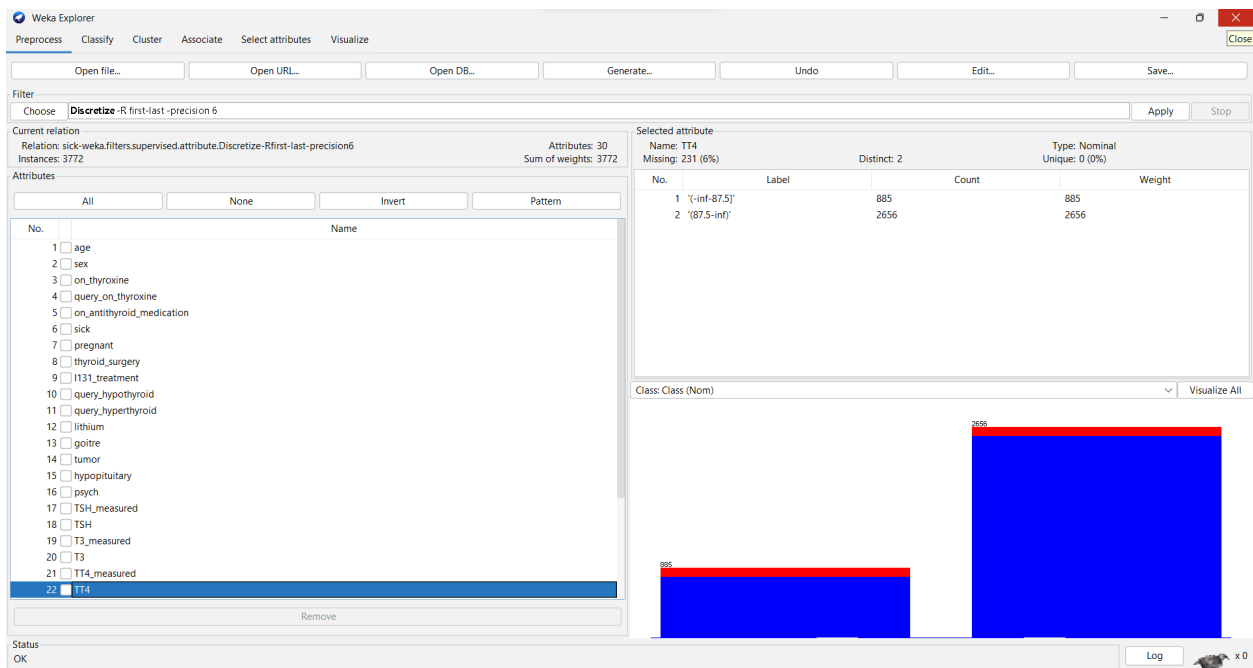
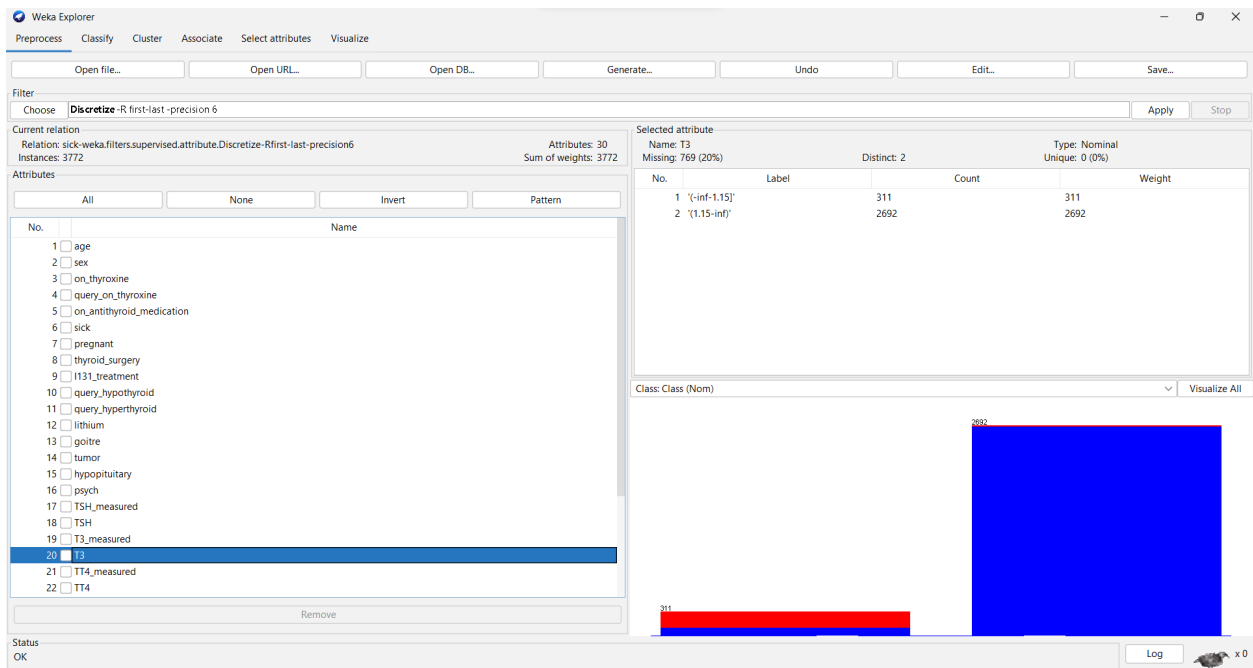
→ sick.arff dataset has been loaded.

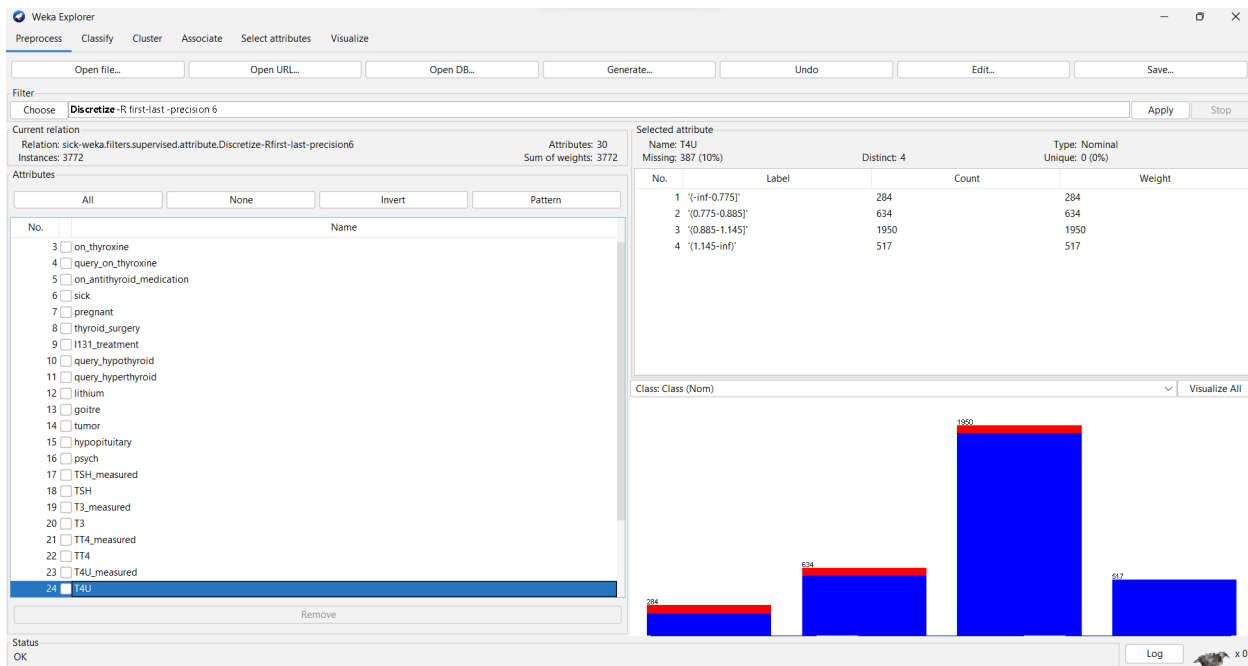


2. Apply the supervised discretization filter on different attributes.

→ Supervised discretization filter has been applied on different attributes. It is used to convert real valued input attributes into nominal attributes. Values are converted from numeric to nominal.







3. What is the effect of this filter on the attributes?

→The Supervised Discretization filter considers the class values and creates distinct boundary ranges for different classes. It is done on the training set and not the test set. It has transformed the numerical variables into categorical counterparts and refers to the target class information. It has simply converted the numeric values into the nominal values. For ex- The filter is applied onto the age attribute where the dataset is converted from numeric to nominal.

4. How many distinct ranges have been created for each attribute?

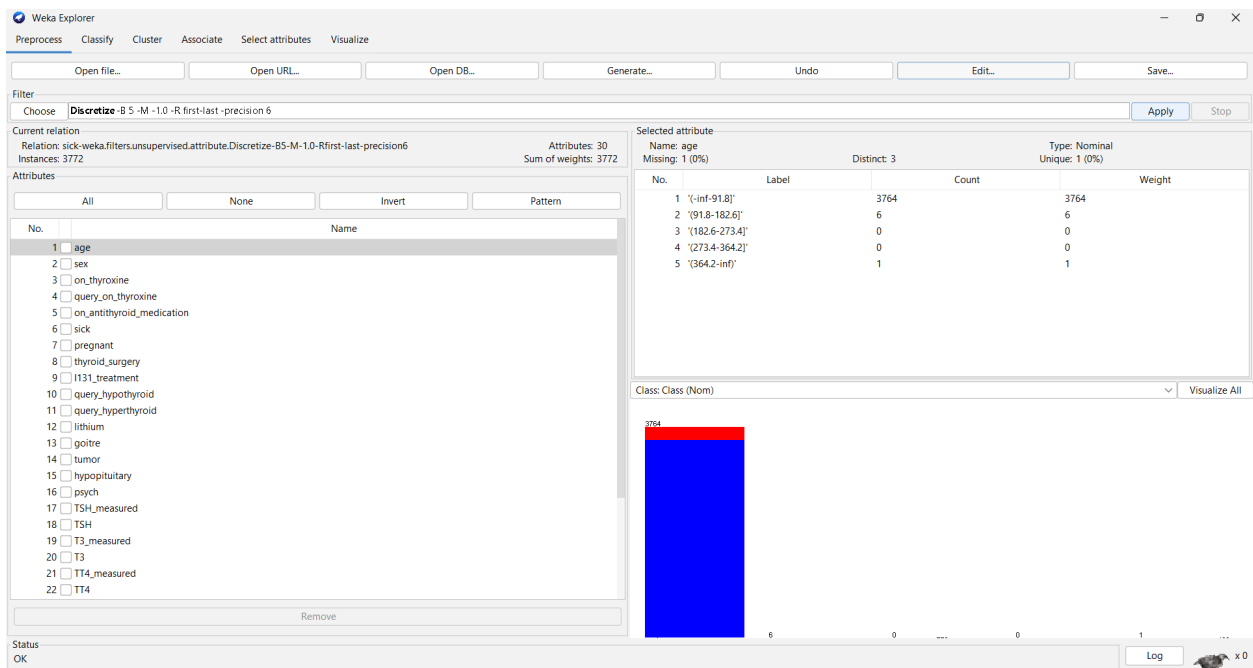
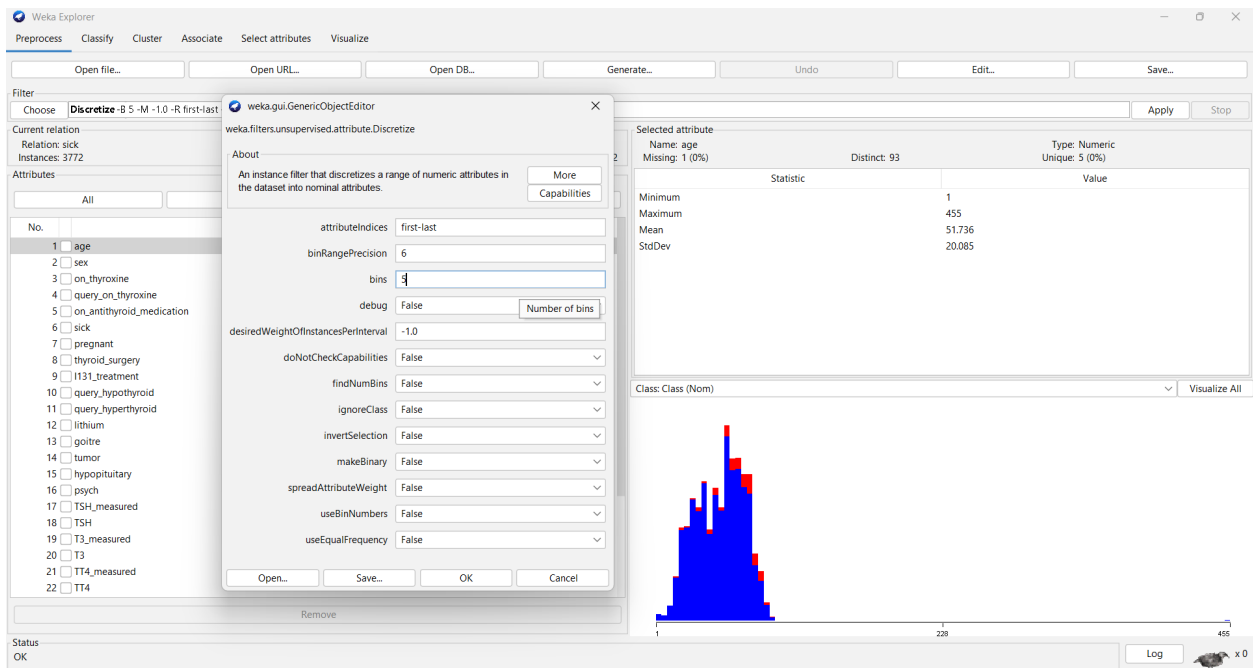
→After Applying the supervised discretization filter distinct classes have been created for each attribute. Consider the age attribute, here three distinct ranges have been created, for TT4 and T3 two distinct ranges with nominal values have been created.

5. Undo the filter applied in the previous step.

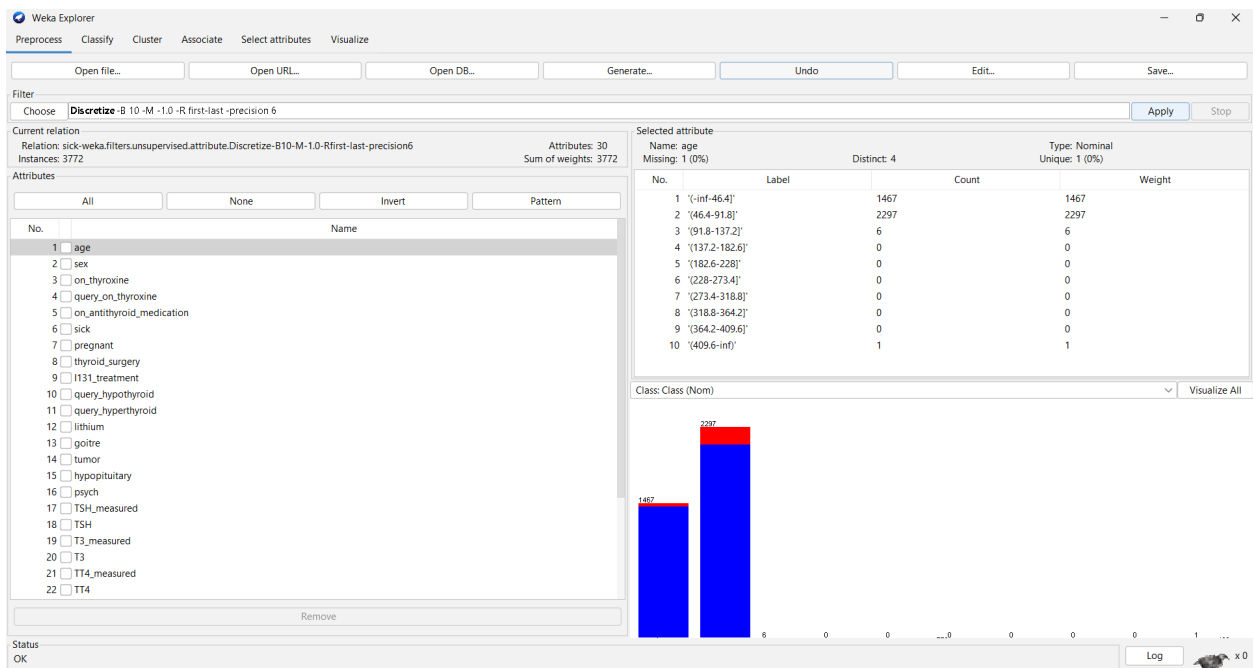
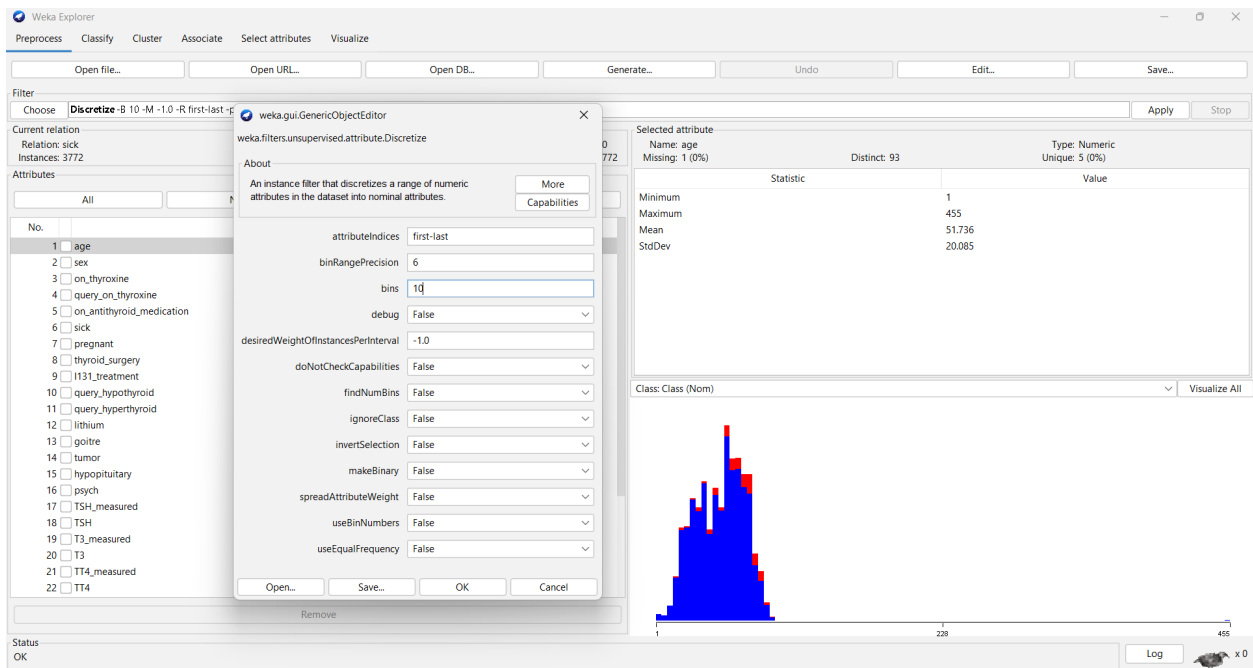
→The UNDO button from the top palette has been pressed.

6. Apply the unsupervised discretization filter. Do this twice:

1. In this step, set 'bins'=5



2. In this step, set 'bins'=10



3. What is the effect of the unsupervised filter on the dataset?

→Unsupervised discretization filters have been applied on the dataset initially with the bin size of 5(effect observed and noted) and then with the bin size of 10 (effect observed on noted). Here the filter has only considered the attribute being discretized. It has not considered the class value.

7. Run the the Naive Bayes classifier after apply the following filters

1. Unsupervised discretized with 'bins'=5

The screenshot shows the Weka Explorer interface with the Naive Bayes classifier selected. The 'Test options' section shows 'Cross-validation' with 'Folds' set to 10. The 'Classifier output' section displays the following results:

```
==== Stratified cross-validation ====
=== Summary ===

Correctly Classified Instances      3455          91.596 %
Incorrectly Classified Instances    317           8.404 %
Kappa statistic                    0.3301
Mean absolute error                 0.1126
Root mean squared error             0.2418
Relative absolute error             97.7 %
Root relative squared error        100.8251 %
Total Number of Instances         3772

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall  F-Measure  MDC   ROC Area  PRC Area  Class
          -----  -
          0.949    0.589    0.961    0.949    0.955    0.332  0.880    0.991    negative
          0.411    0.051    0.344    0.411    0.375    0.332  0.880    0.323    sick
Weighted Avg.   0.916    0.556    0.923    0.916    0.919    0.332  0.880    0.950

=== Confusion Matrix ===

  a    b  <-- classified as
3360 181 |  a = negative
 136   95 |  b = sick
```

The 'Status' bar at the bottom shows 'OK'.

2. Unsupervised discretized with 'bins'=10

The screenshot shows the Weka Explorer interface with the Naive Bayes classifier selected. The 'Test options' section shows 'Cross-validation' with 'Folds' set to 10. The 'Classifier output' section displays the following results:

```
==== Stratified cross-validation ====
=== Summary ===

Correctly Classified Instances      3654          96.8717 %
Incorrectly Classified Instances    118           3.1283 %
Kappa statistic                    0.7405
Mean absolute error                 0.047
Root mean squared error             0.1632
Relative absolute error             40.7549 %
Root relative squared error        68.0853 %
Total Number of Instances         3772

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall  F-Measure  MDC   ROC Area  PRC Area  Class
          -----  -
          0.990    0.203    0.987    0.980    0.983    0.742  0.958    0.997    negative
          0.797    0.020    0.722    0.797    0.757    0.742  0.958    0.677    sick
Weighted Avg.   0.969    0.192    0.970    0.969    0.969    0.742  0.958    0.977

=== Confusion Matrix ===

  a    b  <-- classified as
3470   71 |  a = negative
   47  184 |  b = sick
```

The 'Status' bar at the bottom shows 'OK'.

3. Unsupervised discretized with 'bins'=20

The screenshot shows the Weka Explorer interface with the NaiveBayes classifier selected. The test options are set to Cross-validation with Folds 10 and Percentage split 80. The classifier output shows the following data:

Classifier output		
SVHC	379.0	10.0
other	2169.0	34.0
SVI	849.0	187.0
STMW	113.0	1.0
SVHD	37.0	4.0
[total]	3546.0	236.0

Time taken to build model: 0 seconds

=== Stratified cross-validation ===

=== Summary ===

Metric	Value	Percentage
Correctly Classified Instances	3662	97.0838 %
Incorrectly Classified Instances	110	2.9162 %
Kappa statistic	0.7562	
Mean absolute error	0.0446	
Root mean squared error	0.1596	
Relative absolute error	38.6792 %	
Root relative squared error	66.5739 %	
Total Number of Instances	3772	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
Weighted Avg.	0.982	0.195	0.987	0.982	0.984	0.757	0.965	0.997	negative
	0.805	0.018	0.741	0.805	0.772	0.757	0.965	0.679	sick

=== Confusion Matrix ===

a	b	<-- classified as
3476	65	a = negative
45	106	b = sick

→ For Naive Bayes classification the ratio of training and testing was 80:20. After applying the unsupervised discretization filter the accuracy of the algorithm increased. Accuracy of the algorithm increases with increase in the bin size.

8. Compare the accuracy of the following cases

1. Naive Bayes without discretization filters

The screenshot shows the Weka Explorer interface with the NaiveBayes classifier selected. The test options are set to Cross-validation with Folds 10 and Percentage split 80. The classifier output shows the following data:

Classifier output		
SVHC	379.0	10.0
other	2169.0	34.0
SVI	849.0	187.0
STMW	113.0	1.0
SVHD	37.0	4.0
[total]	3546.0	236.0

Time taken to build model: 0.05 seconds

=== Stratified cross-validation ===

=== Summary ===

Metric	Value	Percentage
Correctly Classified Instances	3493	92.6034 %
Incorrectly Classified Instances	279	7.3966 %
Kappa statistic	0.5249	
Mean absolute error	0.0888	
Root mean squared error	0.2254	
Relative absolute error	77.0863 %	
Root relative squared error	95.6866 %	
Total Number of Instances	3772	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
Weighted Avg.	0.936	0.225	0.905	0.936	0.960	0.550	0.925	0.991	negative
	0.775	0.064	0.441	0.775	0.562	0.550	0.925	0.660	sick

=== Confusion Matrix ===

a	b	<-- classified as
3314	227	a = negative
52	179	b = sick

2. Naive Bayes with a supervised discretization filter

The screenshot shows the Weka Explorer interface with the Naive Bayes classifier selected. The 'Test options' section shows 'Cross-validation' with 'Folds' set to 10. The 'Classifier output' section displays the following results:

```
==== Stratified cross-validation ====
==== Summary ====
Correctly Classified Instances      3670           97.2959 %
Incorrectly Classified Instances      102           2.7041 %
Kappa statistic                    0.7748
Mean absolute error                  0.0439
Root mean squared error              0.1574
Relative absolute error              38.069 %
Root relative squared error          65.6429 %
Total Number of Instances          3772

==== Detailed Accuracy By Class ====
              TP Rate  FP Rate  Precision  Recall  F-Measure  MDC   ROC Area  PRC Area  Class
negative      0.982    0.173    0.989    0.982    0.986    0.776    0.960    0.997
sick           0.827    0.018    0.755    0.827    0.789    0.776    0.960    0.733
Weighted Avg.  0.973    0.164    0.974    0.973    0.974    0.776    0.960    0.980

==== Confusion Matrix ====
a  b  <-- classified as
3479  62 | a = negative
  40 191 | b = sick
```

The status bar at the bottom shows 'OK'.

3. Naive Bayes with an unsupervised discretization filter with different values for the bins attributes.

1. Unsupervised discretized with 'bins'=5

The screenshot shows the Weka Explorer interface with the Naive Bayes classifier selected. The 'Test options' section shows 'Cross-validation' with 'Folds' set to 10. The 'Classifier output' section displays the following results:

```
==== Stratified cross-validation ====
==== Summary ====
Correctly Classified Instances      3455           91.596 %
Incorrectly Classified Instances      317           8.404 %
Kappa statistic                    0.3301
Mean absolute error                  0.1126
Root mean squared error              0.2418
Relative absolute error              97.7 %
Root relative squared error          100.8251 %
Total Number of Instances          3772

==== Detailed Accuracy By Class ====
              TP Rate  FP Rate  Precision  Recall  F-Measure  MDC   ROC Area  PRC Area  Class
negative      0.949    0.589    0.961    0.949    0.955    0.332    0.880    0.991
sick           0.411    0.051    0.344    0.411    0.375    0.332    0.880    0.323
Weighted Avg.  0.916    0.556    0.923    0.916    0.919    0.332    0.880    0.950

==== Confusion Matrix ====
a  b  <-- classified as
3360 181 | a = negative
  136  95 | b = sick
```

The status bar at the bottom shows 'OK'.

2. Unsupervised discretized with 'bins'=10

The screenshot shows the Weka Explorer interface with the NaiveBayes classifier selected. The 'Test options' section on the left has 'Cross-validation' selected with 'Folds' set to 10. The 'Classifier output' pane on the right displays the following results:

Time taken to build model: 0 seconds

=== Stratified cross-validation ===

=== Summary ===

Metric	Value
Correctly Classified Instances	3654
Incorrectly Classified Instances	118
Kappa statistic	0.7405
Mean absolute error	0.047
Root mean squared error	0.1632
Relative absolute error	40.7549 %
Root relative squared error	65.0853 %
Total Number of Instances	3772

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
Weighted Avg.	0.980	0.203	0.987	0.980	0.983	0.742	0.958	0.997	negative
	0.797	0.020	0.722	0.797	0.757	0.742	0.958	0.677	sick

=== Confusion Matrix ===

	a	b	<-- classified as
3470	71		a = negative
47	184		b = sick

3. Unsupervised discretized with 'bins'=20

The screenshot shows the Weka Explorer interface with the NaiveBayes classifier selected. The 'Test options' section on the left has 'Cross-validation' selected with 'Folds' set to 10. The 'Classifier output' pane on the right displays the following results:

Time taken to build model: 0 seconds

=== Stratified cross-validation ===

=== Summary ===

Metric	Value
Correctly Classified Instances	3662
Incorrectly Classified Instances	110
Kappa statistic	0.7562
Mean absolute error	0.0446
Root mean squared error	0.1556
Relative absolute error	38.6792 %
Root relative squared error	66.5739 %
Total Number of Instances	3772

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
Weighted Avg.	0.982	0.195	0.987	0.982	0.984	0.757	0.965	0.997	negative
	0.805	0.018	0.741	0.805	0.772	0.757	0.965	0.679	sick

=== Confusion Matrix ===

	a	b	<-- classified as
3476	65		a = negative
45	186		b = sick

→

Done in question 7 with bin size 5,10,20.

Accuracy in Naive bayes without discretisation is lesser than that with supervised discretization. Accuracy further increased with unsupervised discretization than that of supervised as the bin number increased.

Accuracy of Naive Bayes without discretization filters is: 92.6034 %

Accuracy of Naive Bayes with a supervised discretization filter: 97.2959 %

Accuracy of Naive Bayes with an unsupervised discretization filter with different values for the bins attributes increases with increase in bin size. with bins=5, accuracy is 91%(approx) and with bin size=20, accuracy is 97% (approx).
