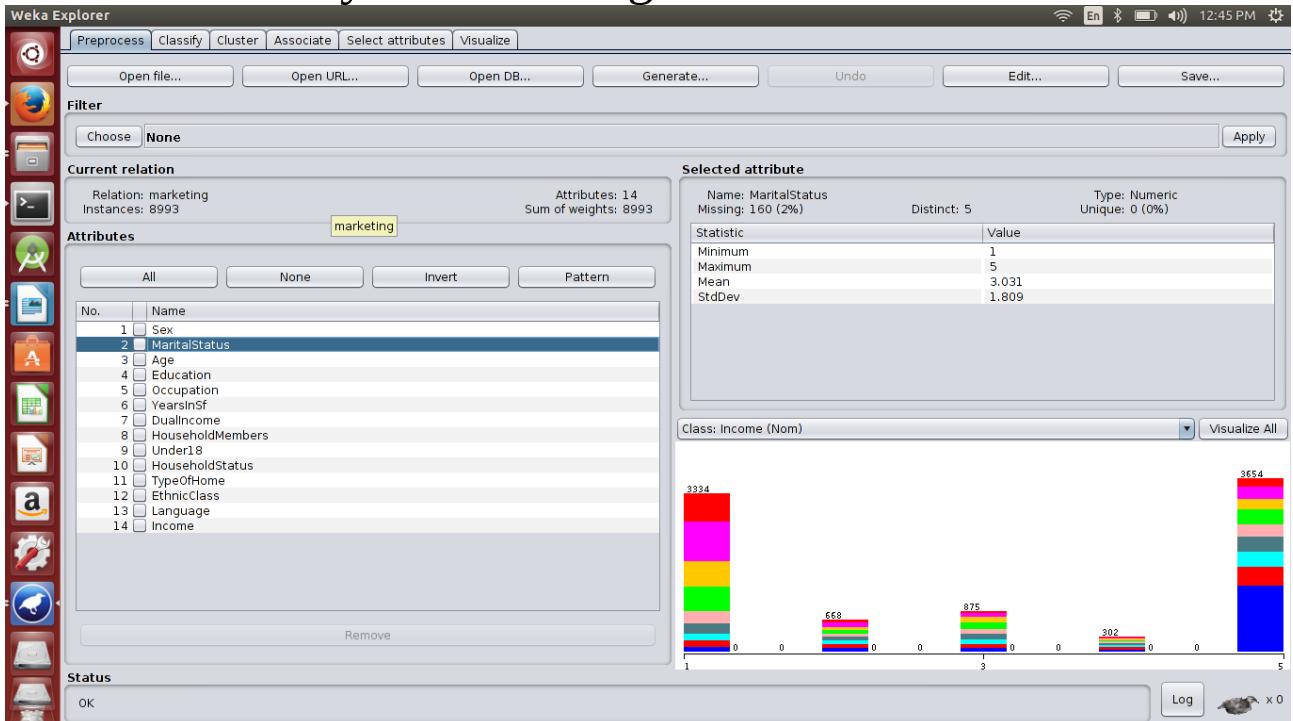


# DM ASSIGNMENT

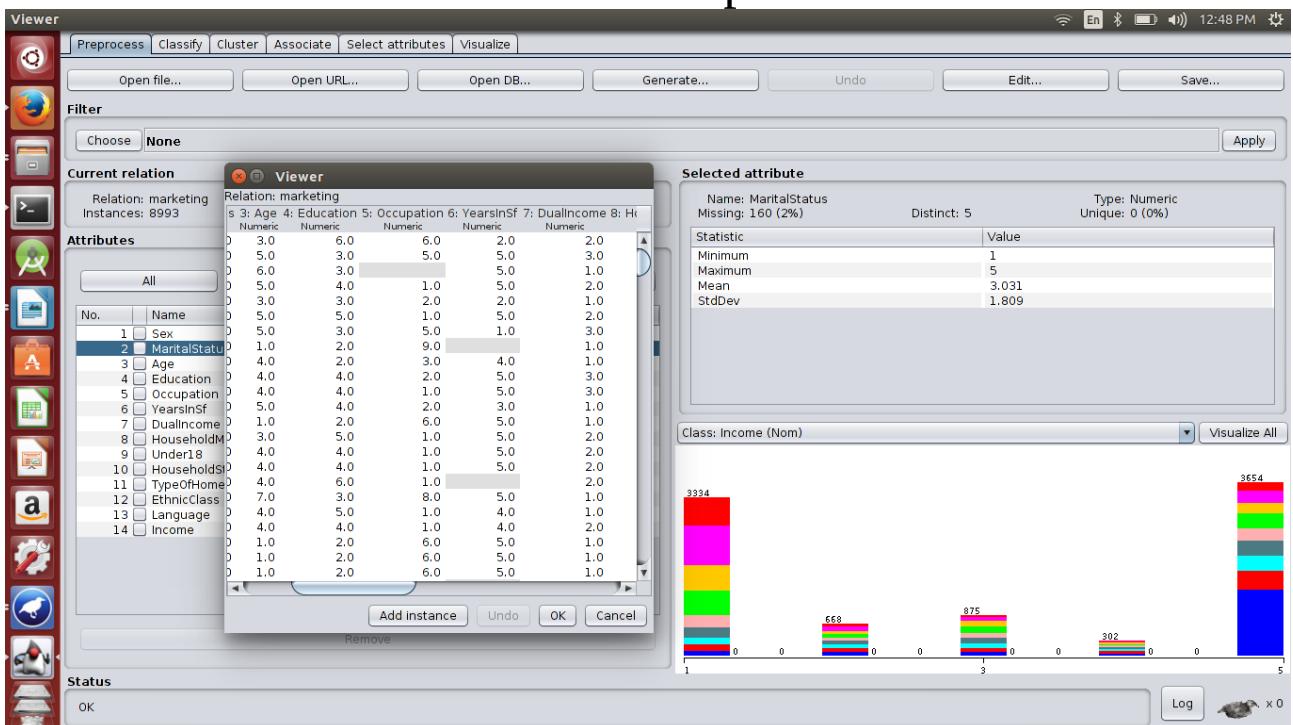
Question-3:

1) Preprocess missing values with the various options provided by WEKA:

Data Set initially with missing values:



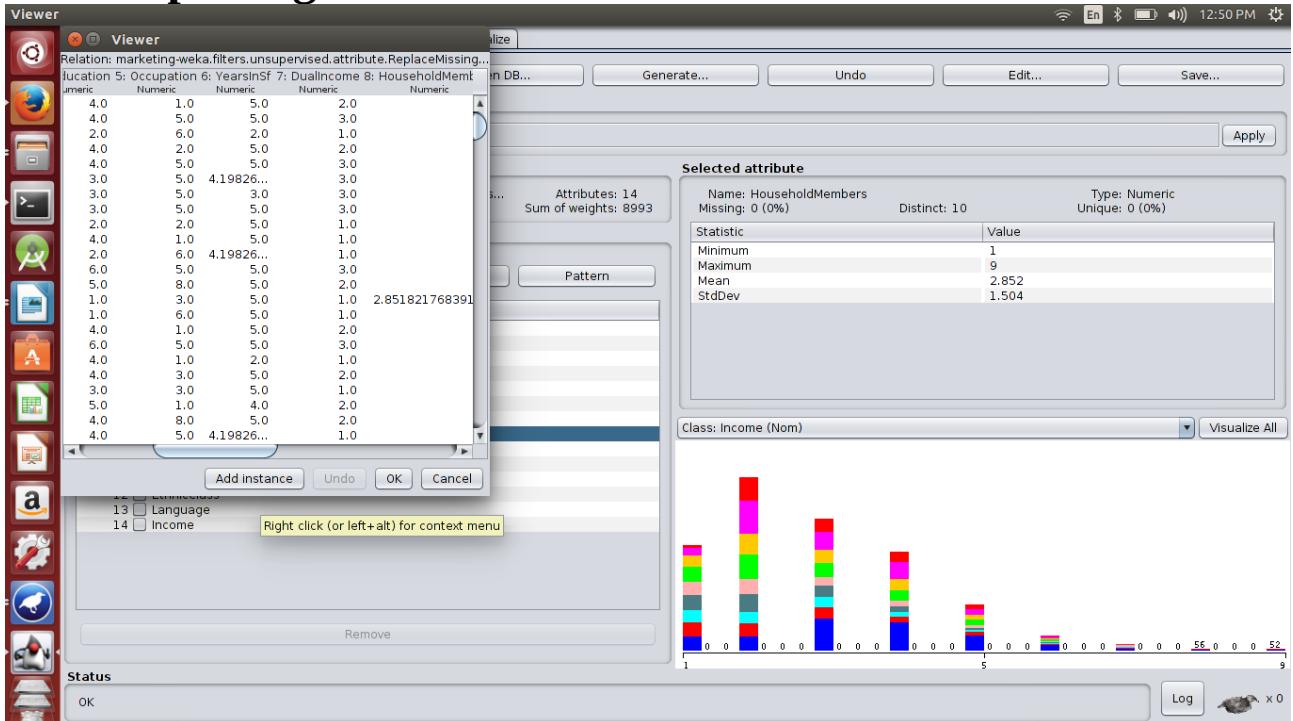
Use the Edit button to view the tuples.



We can do this using ReplaceMissingValues filter. Click the

“Choose” button for the Filter in preprocess tab and select ReplaceMissingValues, it us under unsupervised.attribute.ReplaceMissingValues.

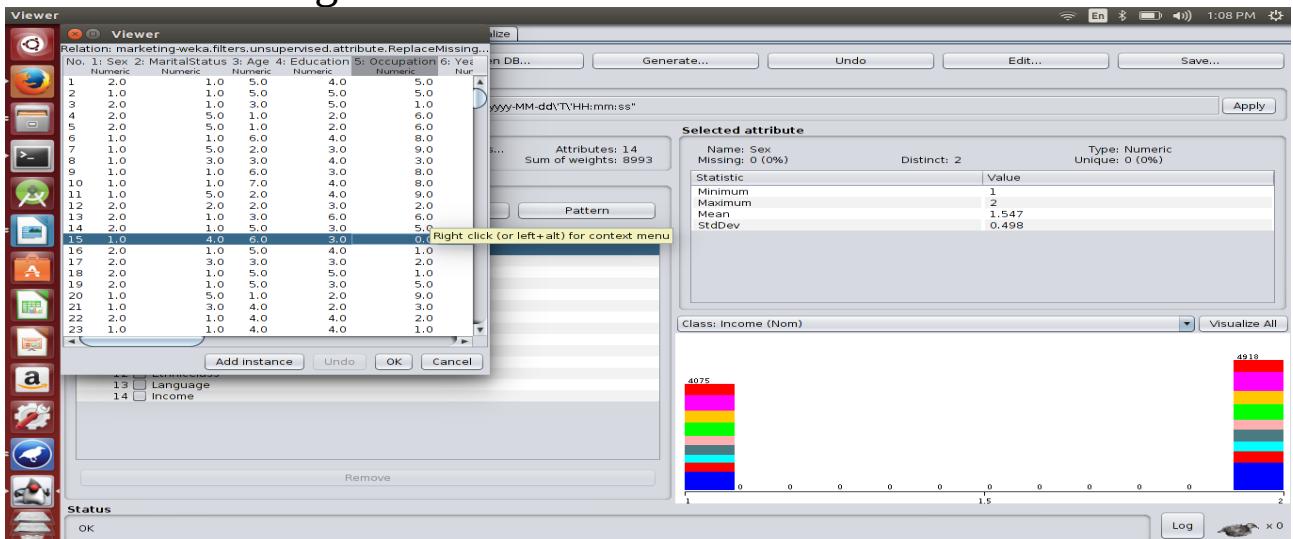
**After replacing:**



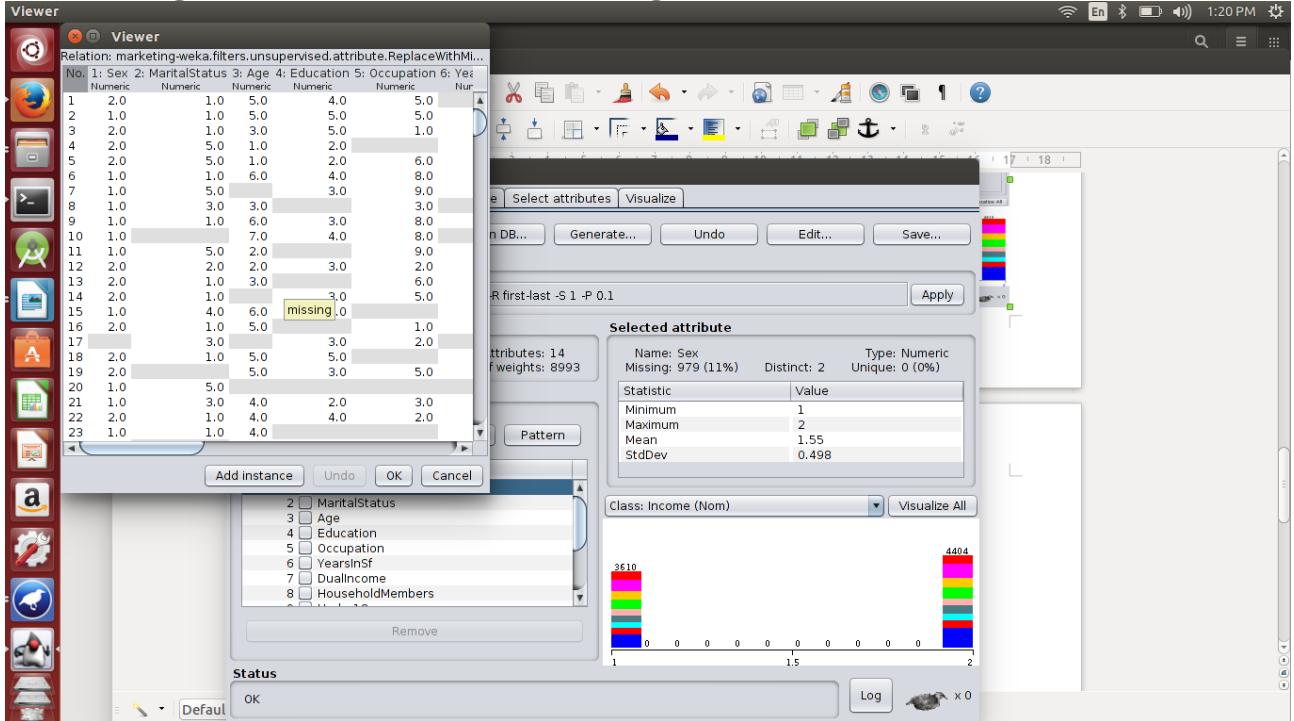
So, the values got filled.

### Using ReplaceMissingwith User Constant filter:

Click the “Choose” button for the Filter in preprocess tab and select ReplaceMissingwith User Constant it us under unsupervised.attribute.ReplaceMissingwith User Constant . This fills missing values with 0.0 in this case

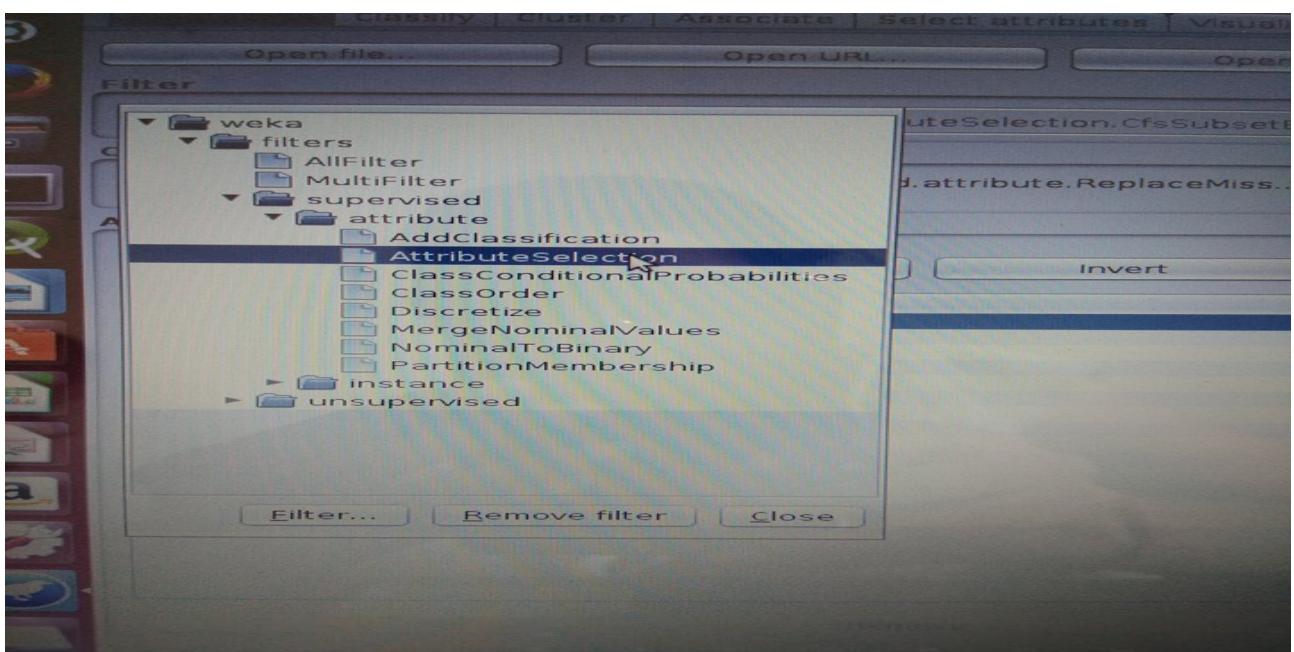


Where as using replace with missing values does the following. Increases no: of missing values.



## 2)Attribute Filter option of WEKA :

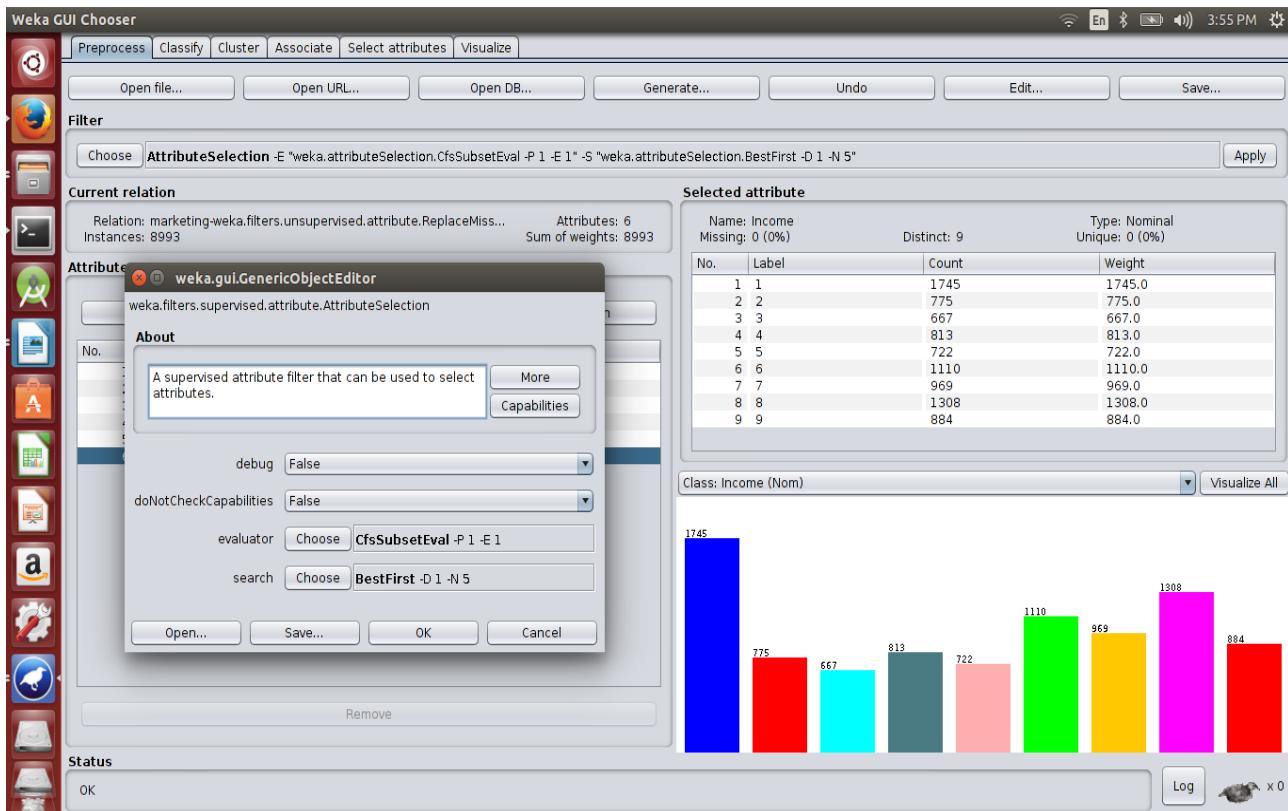
Select Attribute selection option from attribute which is in supervised filter in the preprocess tab.



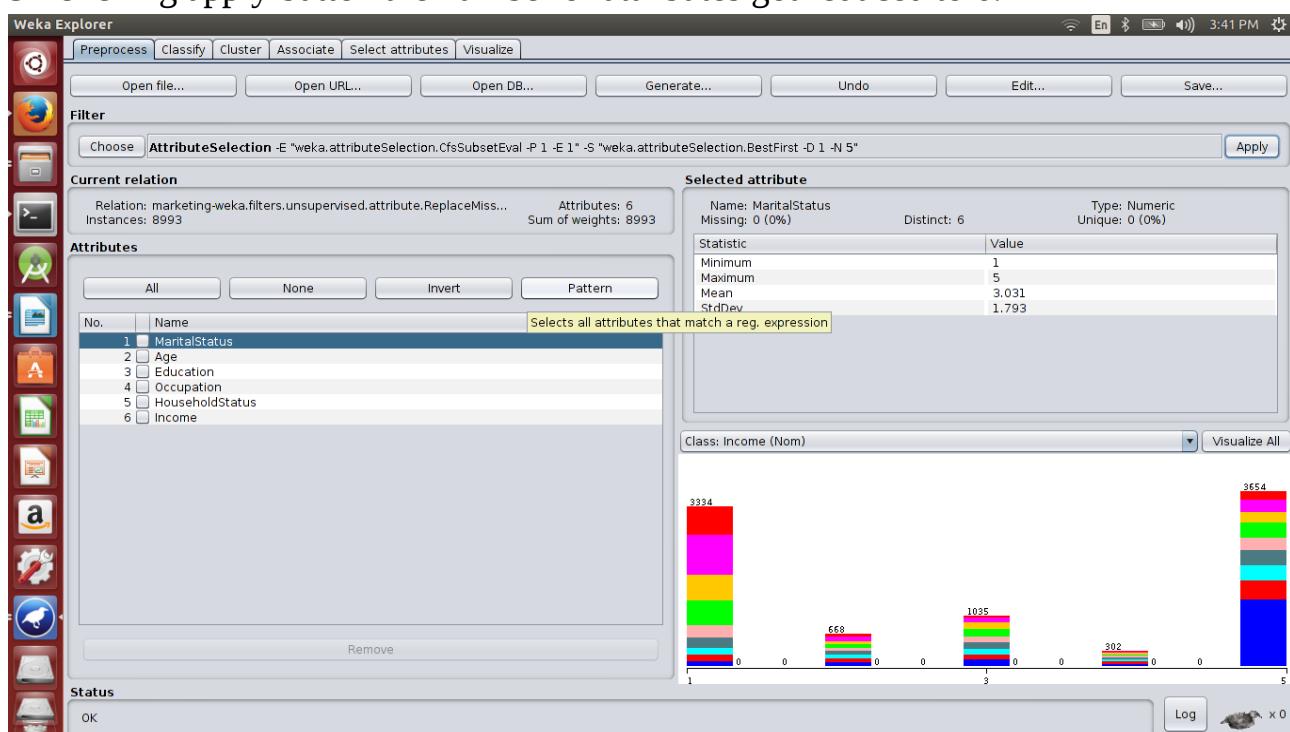
This default uses CfsSubsetEval and BestFirst as Search.

## CfsSubsetEval :

Evaluates the worth of a subset of attributes by considering the individual predictive ability of each feature along with the degree of redundancy between them.



On clicking apply button the number of attributes got reduced to 6.

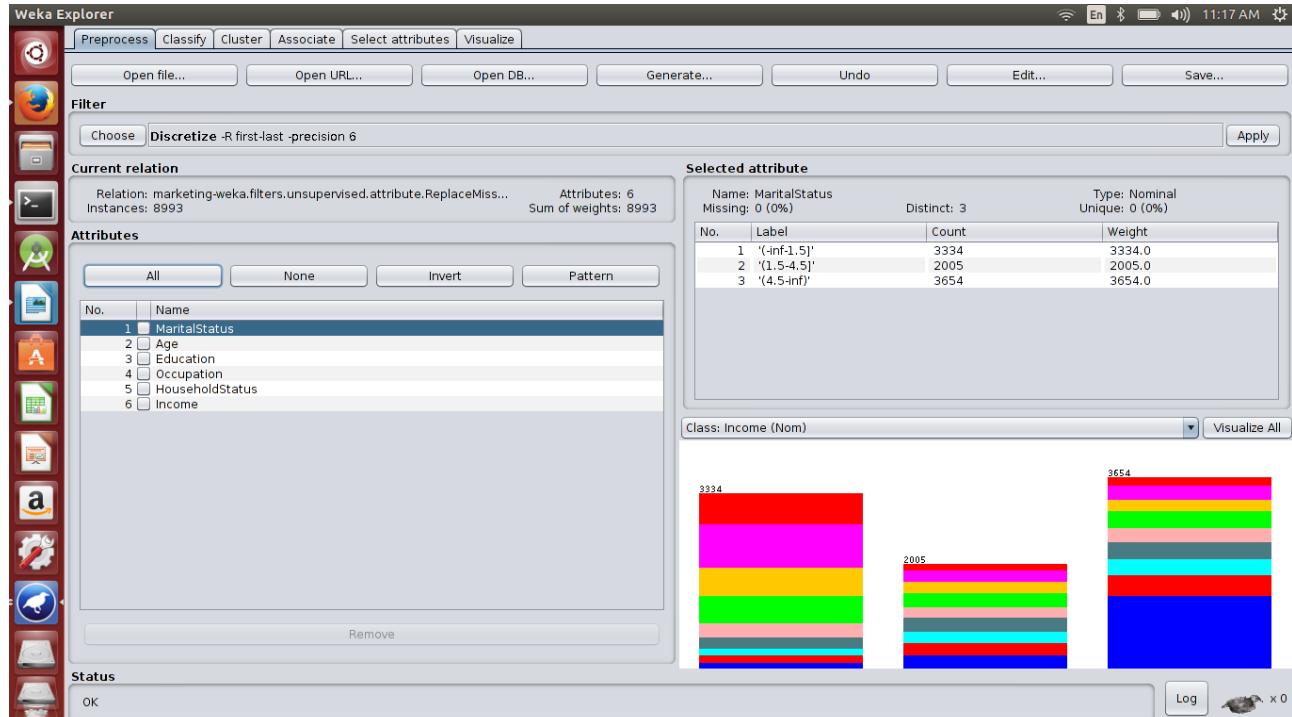


### 3)Discretization in Weka:

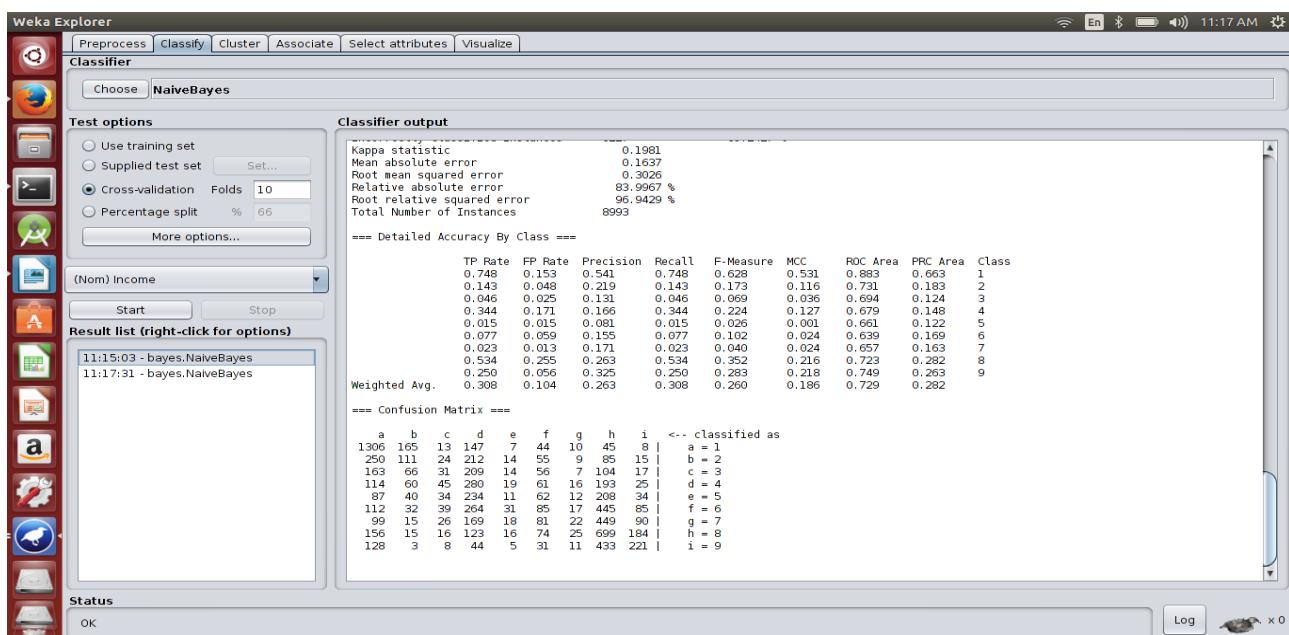
One has numeric data but wants to use classifier that handles only nominal values. In that case one needs to *discretize* the data. This can be done in two ways:

#### 1)weka.filters.supervised.attribute.Discretize

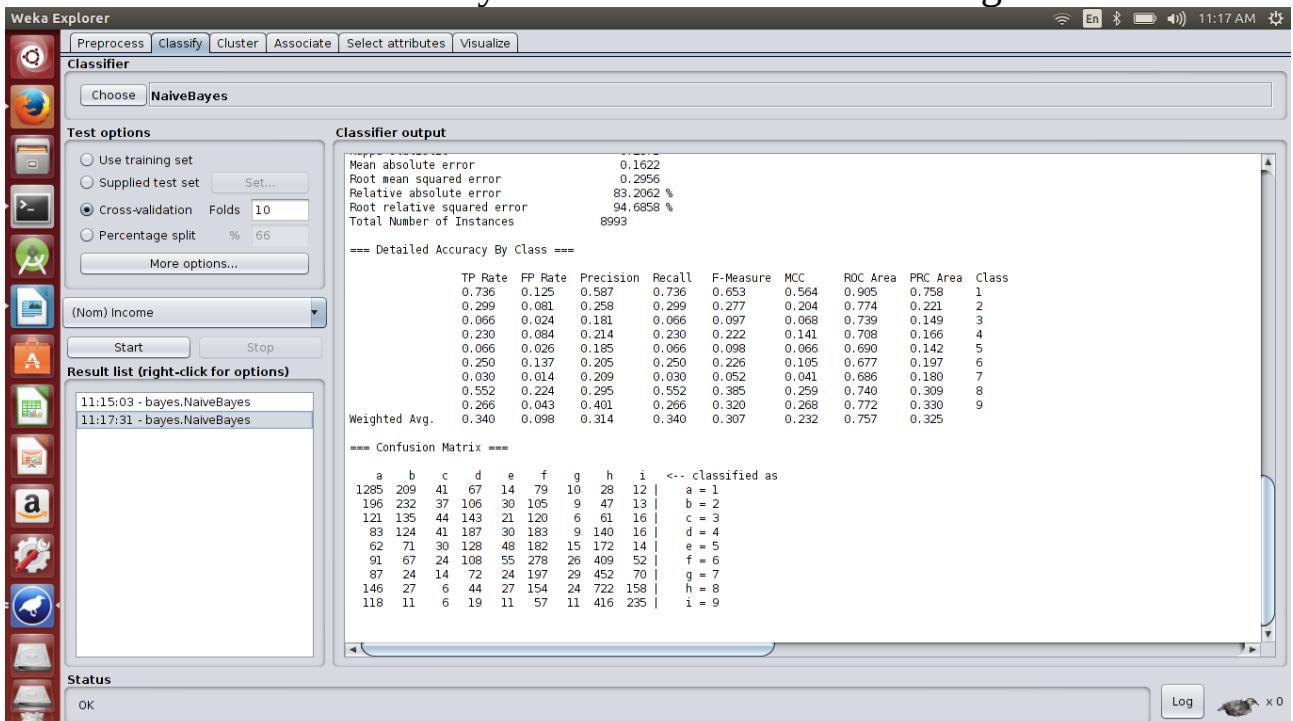
uses either Fayyad & Irani's MDL method or Kononeko's MDL criterion. This supervised Discretize option considers the correlation between attributes and class Attribute.



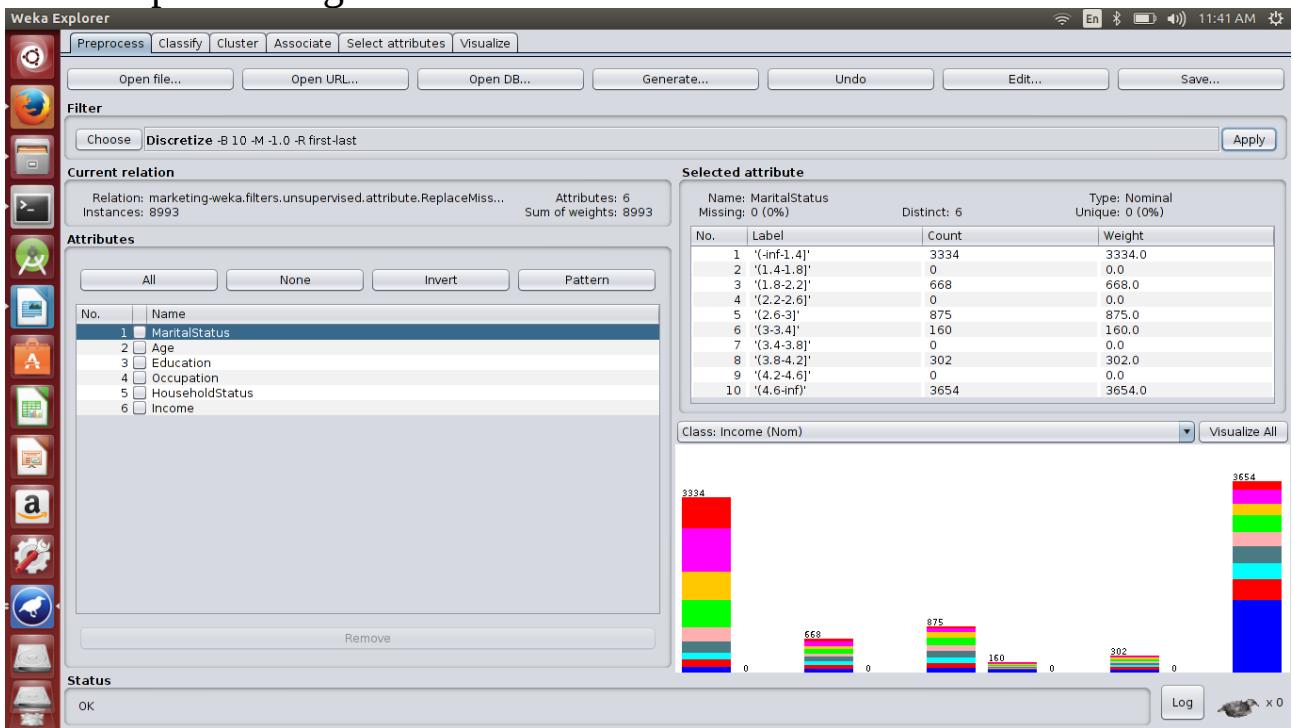
So, on applying discretize values of certain ranges are classified under same category as seen in above image. Ex'('(-inf,1.5)'



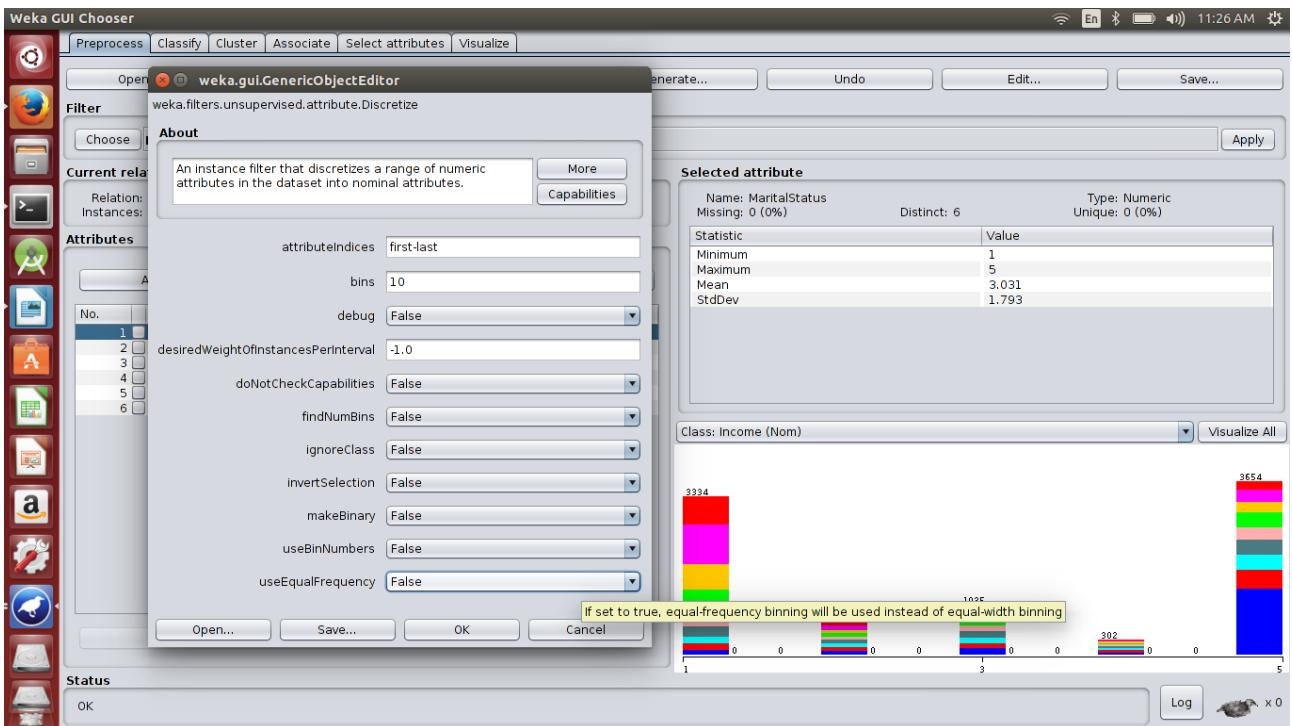
Before discretization F-measure is 0.260 which increased to 0.307 after discretize which clearly shows the use of discretizing values.



2) weka.filters.unsupervised.attribute.Discretize  
uses simple binning

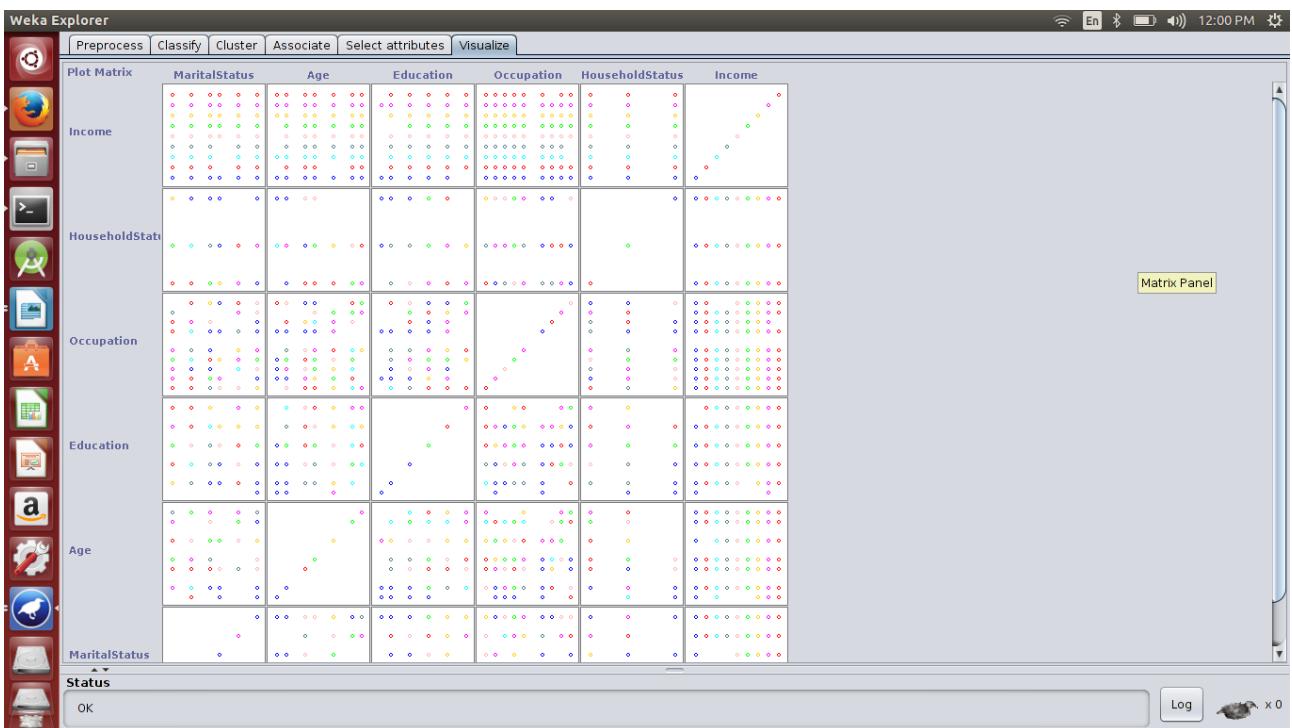


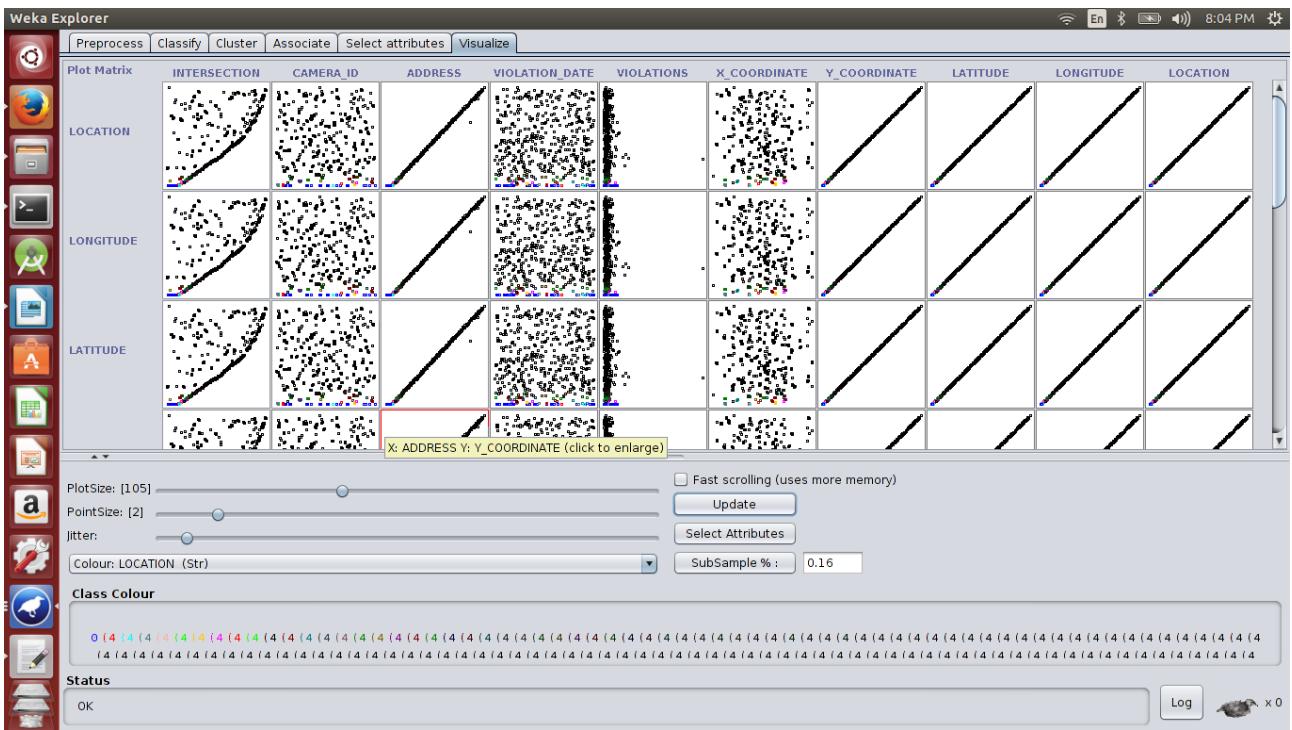
We have options to select no: of bins and to use equal frequency or equal width. This on the other side does not consider correlation between attributes and class Attribute.



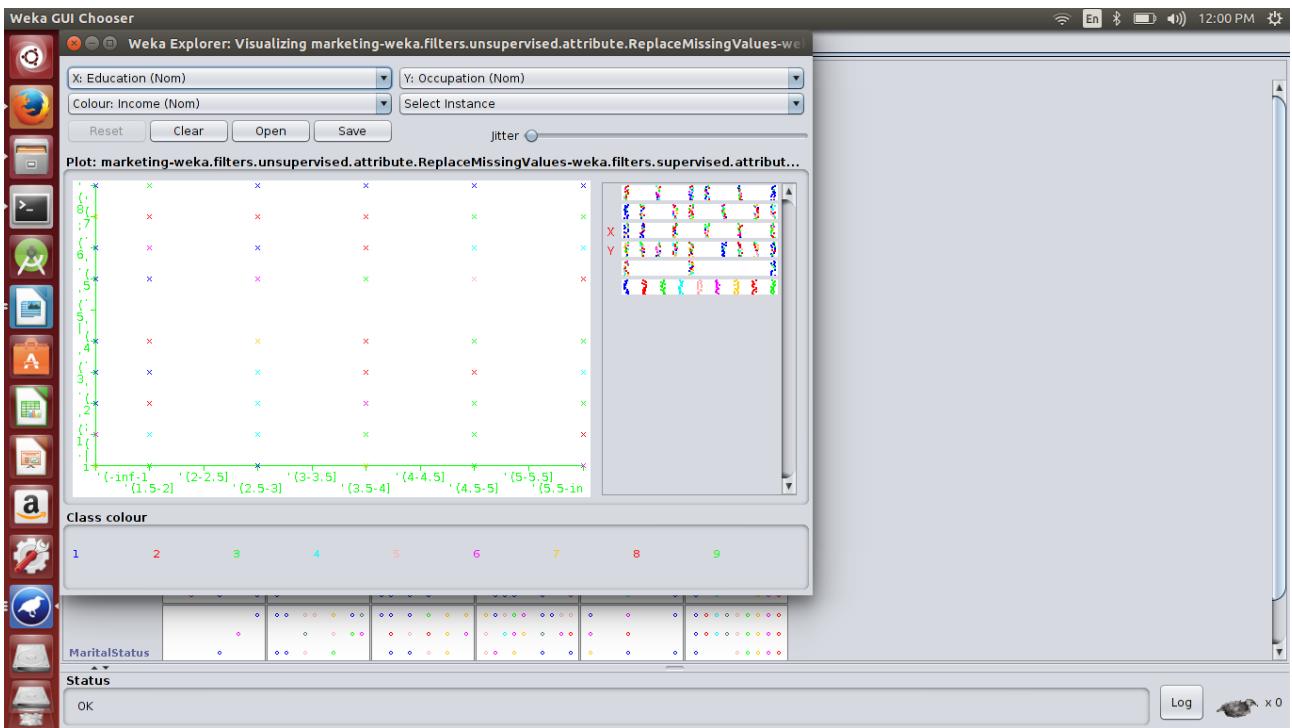
#### 4) Various Visualization techniques supported by WEKA :

Can highlight any patterns in the relationship between the attributes, such as positive or negative correlations. So, Weka provides us Visualize tab for this purpose.





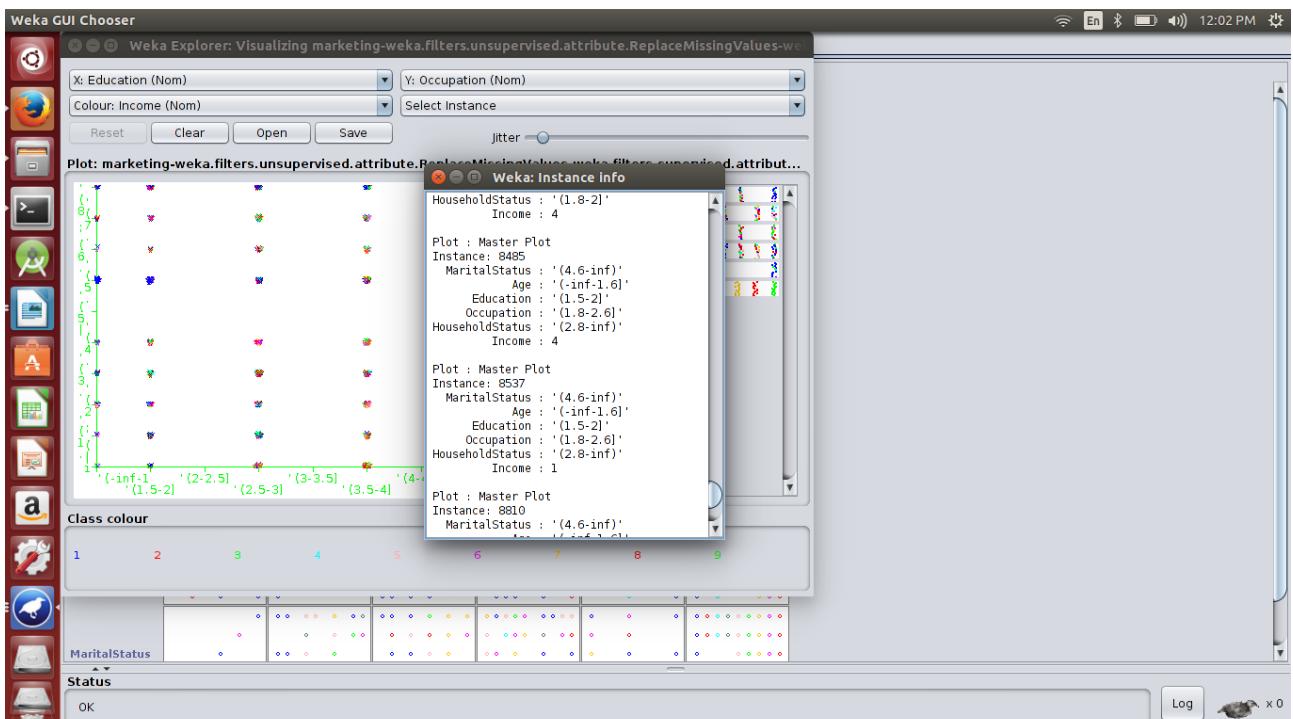
The dots in the plots are colored by their class value.  
 Clicking on a plot will give you a new window with the plot.  
 all combinations of attributes are plotted in a systematic way.



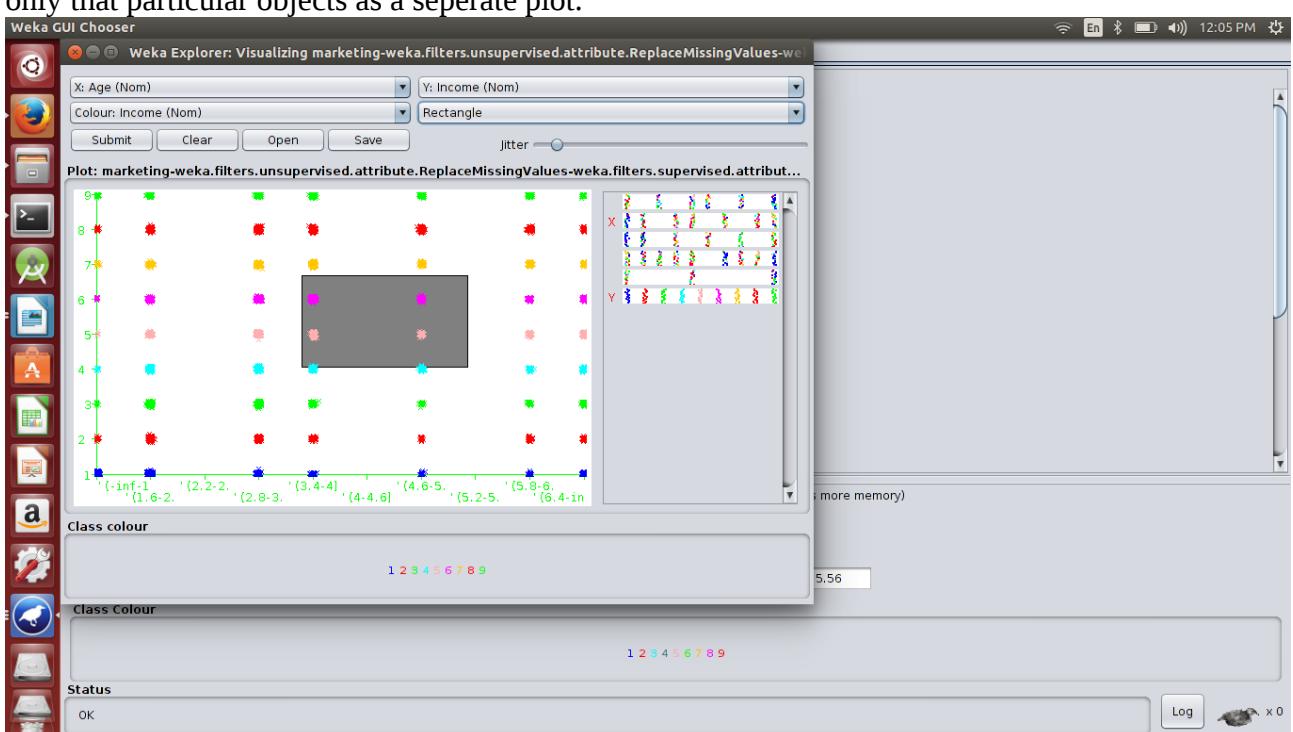
We can also see that each plot appears twice, first in the top left triangle and again in the bottom right triangle with the axes flipped. We can also see a series of plots starting in the bottom left and continuing to the top right where each attribute is plotted

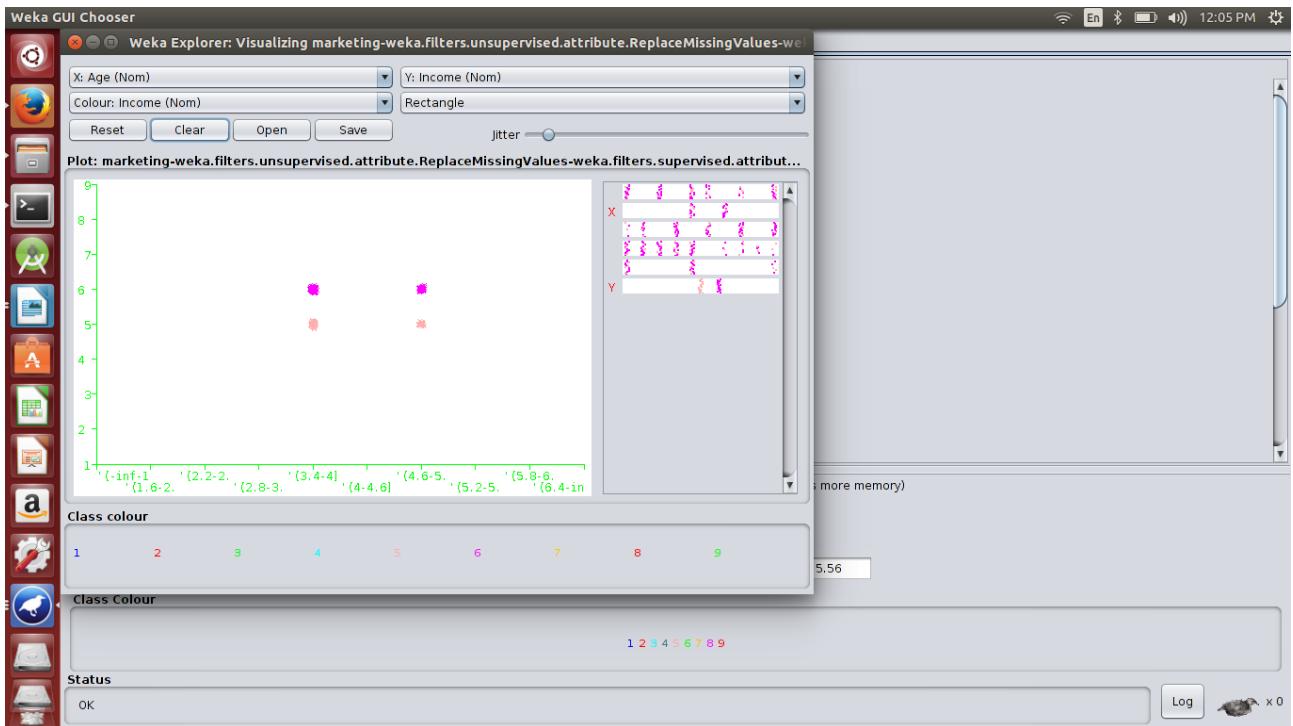
against itself.

Controls at the bottom of the screen. These let us increase the size of the plots, increase the size of the dots and add jitter. On increasing jitter and clicking on a particular point we get all the details of instances which are having same particular x and y values.

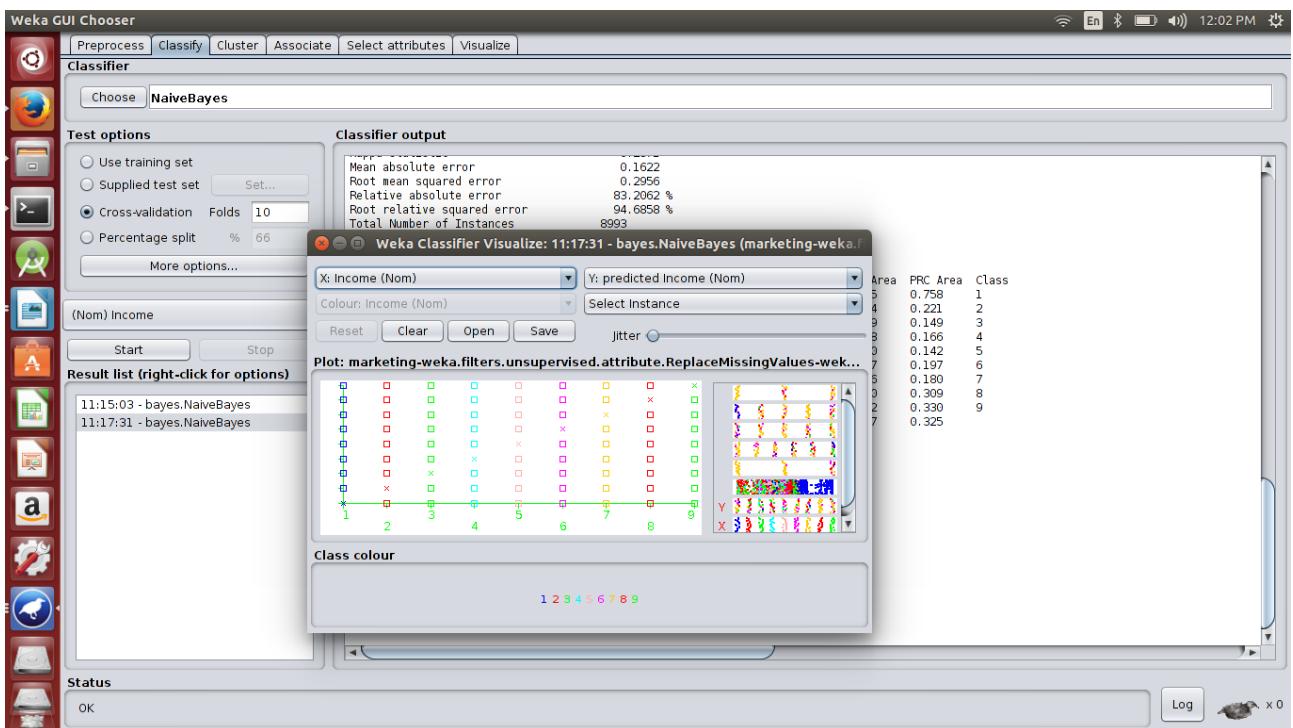


We can select the instance by choosing a shape ex:rectangle so on clicking submit we get values of only that particular objects as a separate plot.



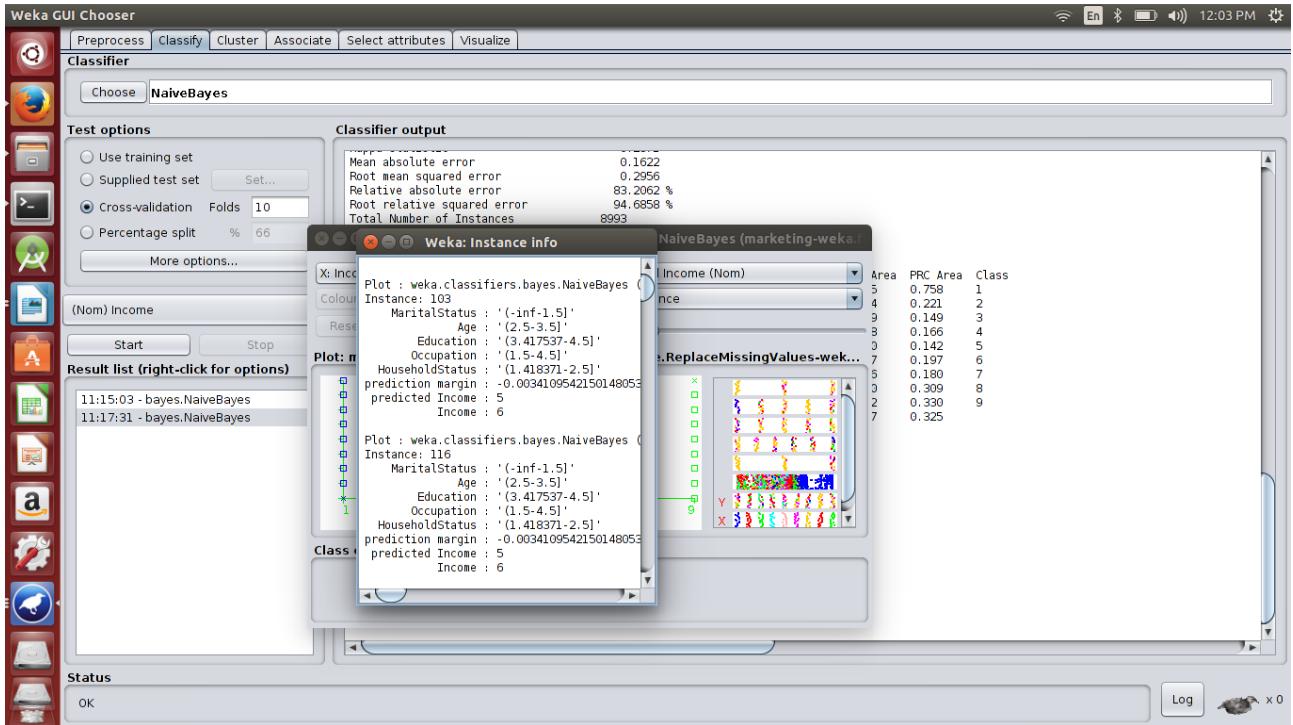


**We can Visualize classification errors so that which helps in proper data exploring and so that we get to know which all instances are classified wrong and we can correct them.**  
**So right click on the classifier then select Visualize Classification Errors, and we obtain the following plot:**

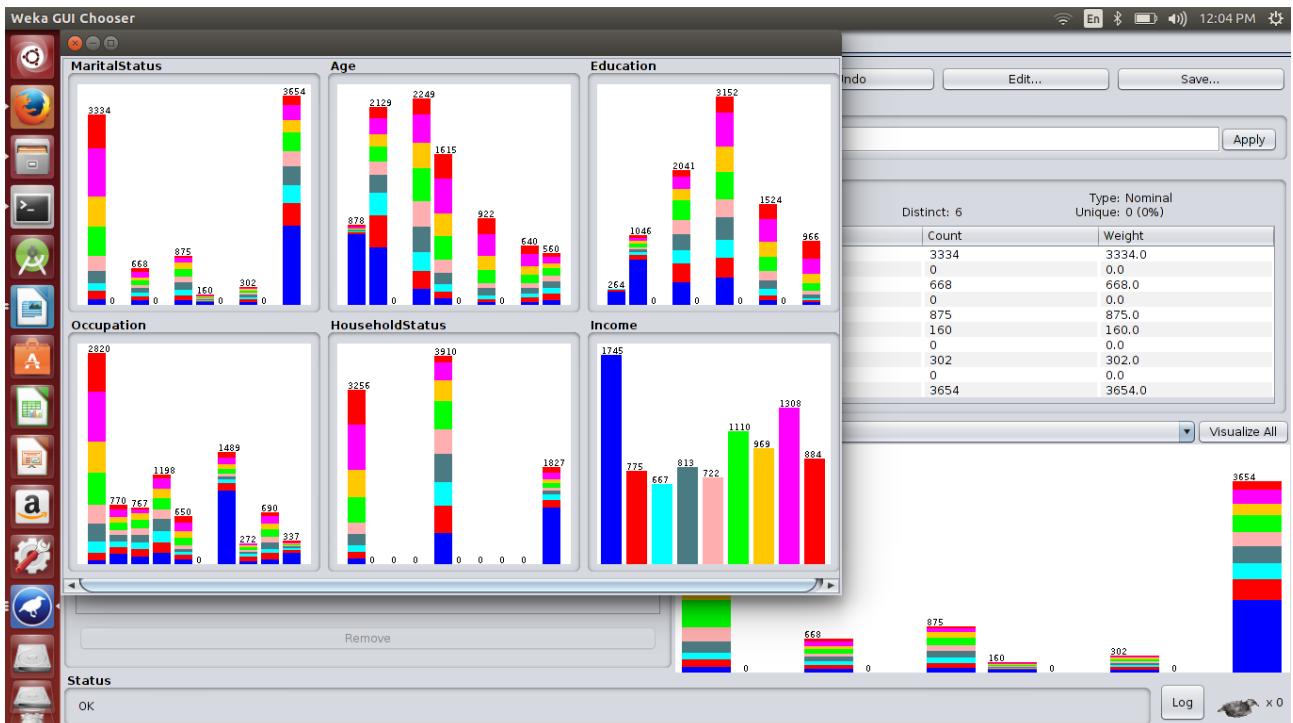


On clicking on the point we get to know what actual class it should be and what

actually it is being predicted.



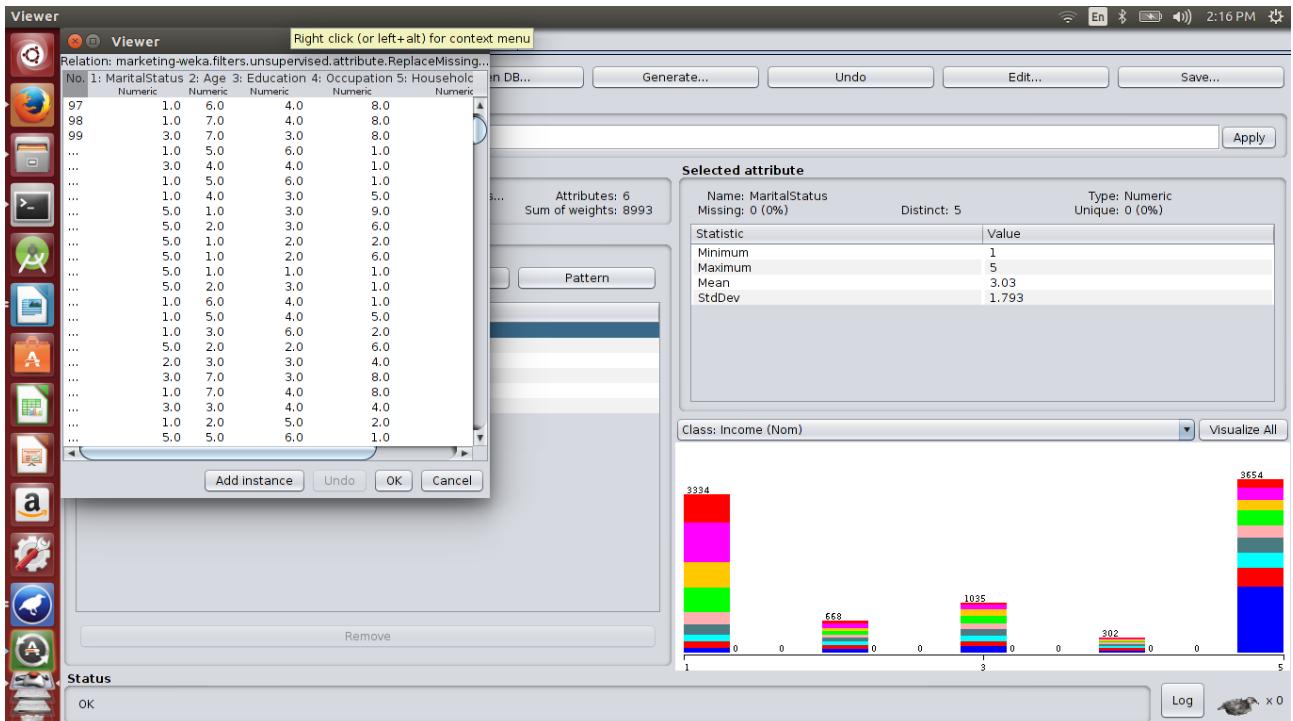
In the preprocess tab on clicking Visualize all we get the histogram of each attribute and their class values are indicated by colours.



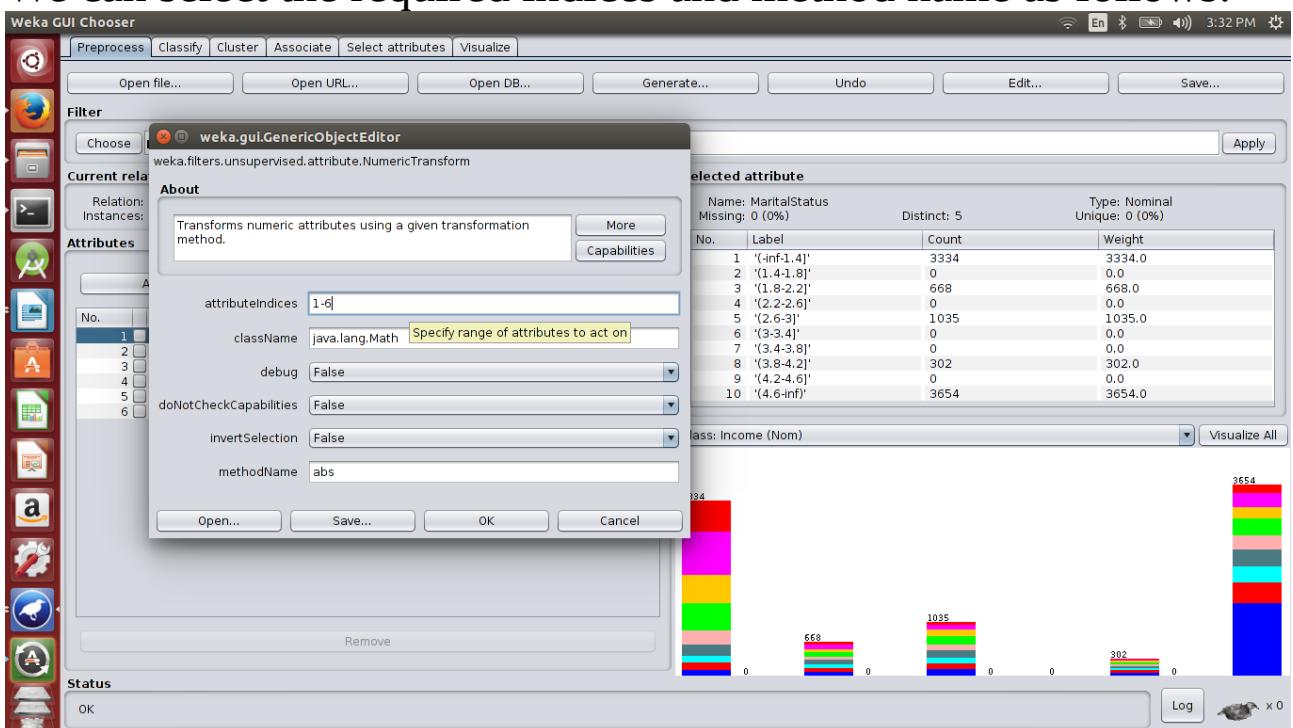
## 5) Data Transformation by WEKA:

As the given data consists all numeric attributes and as we filled the missing values with mean of all other values ,so as the average will be in decimal in case if we need it to convert to absolute value or floor of the value we use numeric transform.We can select this from unsupervised.attribute.numeric transform

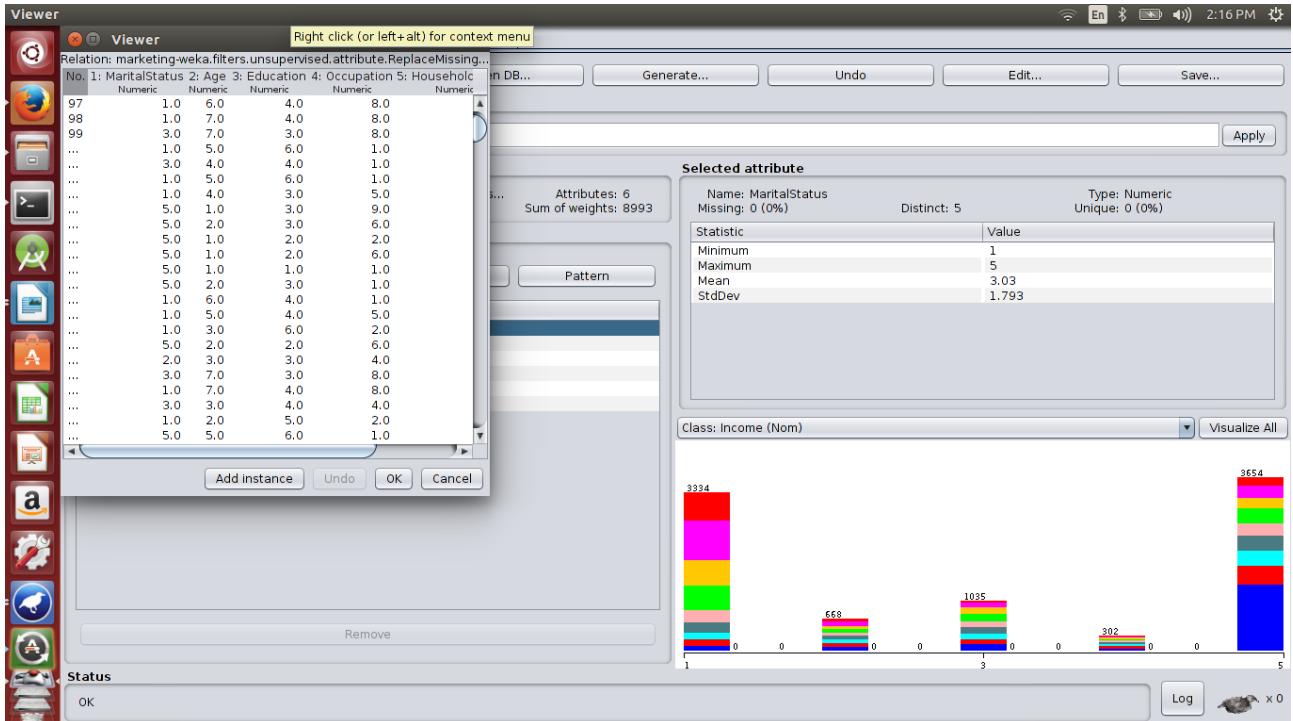
Before numeric transform:



We can select the required indices and method name as follows:

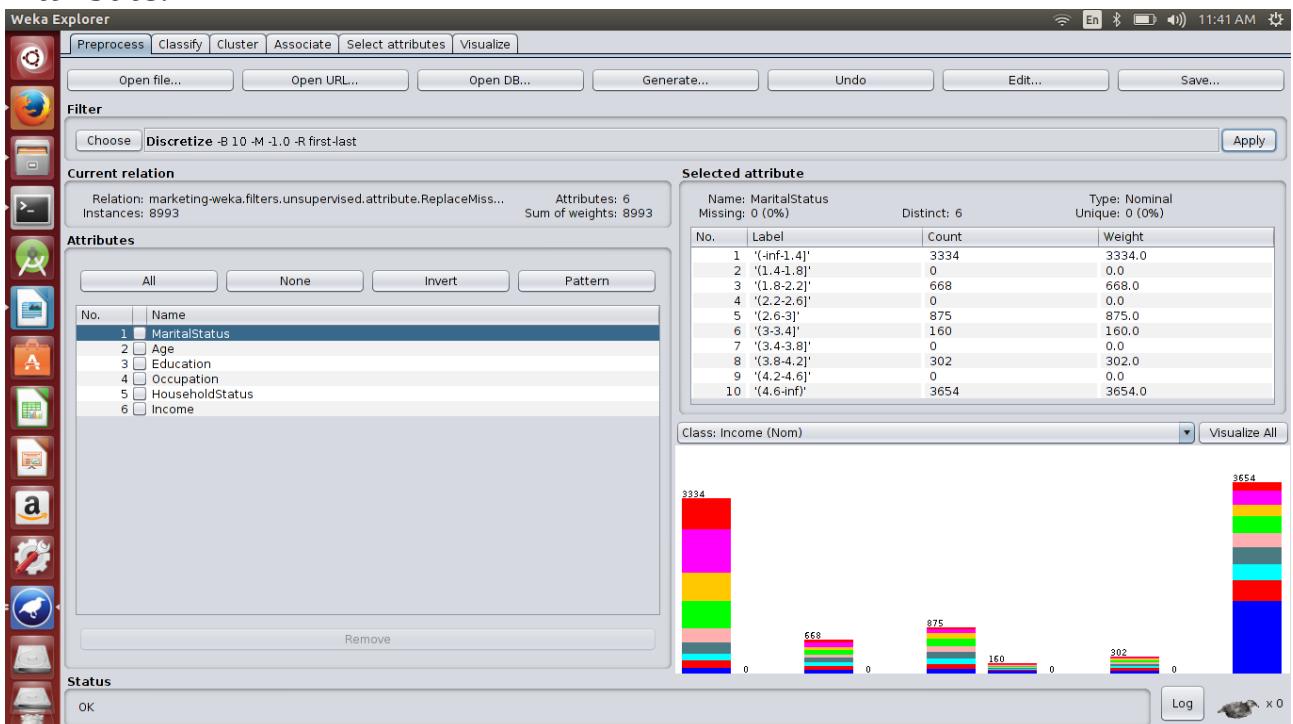


## After numeric transform:



Also discretization comes under Data transform

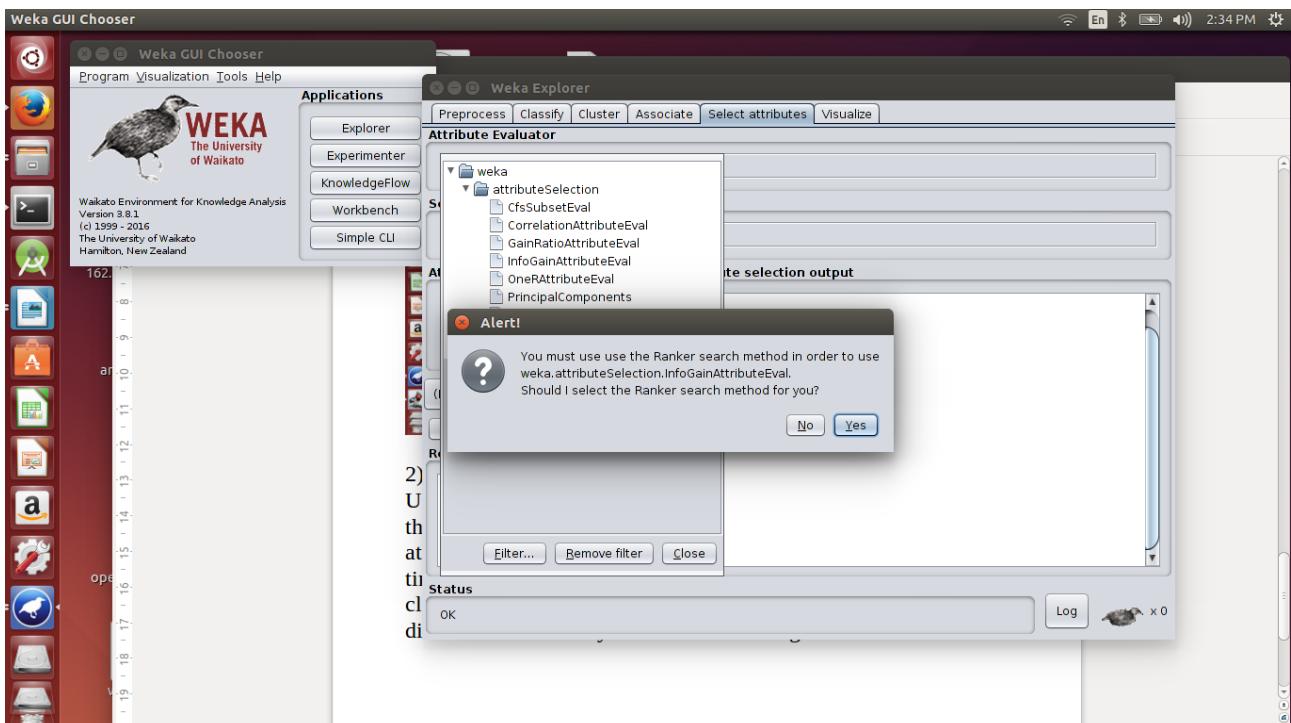
`weka.filters.unsupervised.attribute.Discretize` uses simple binning. We have options to select no: of bins and to use equal frequency or equal width. This on the other side does not consider correlation between attributes and class Attribute.



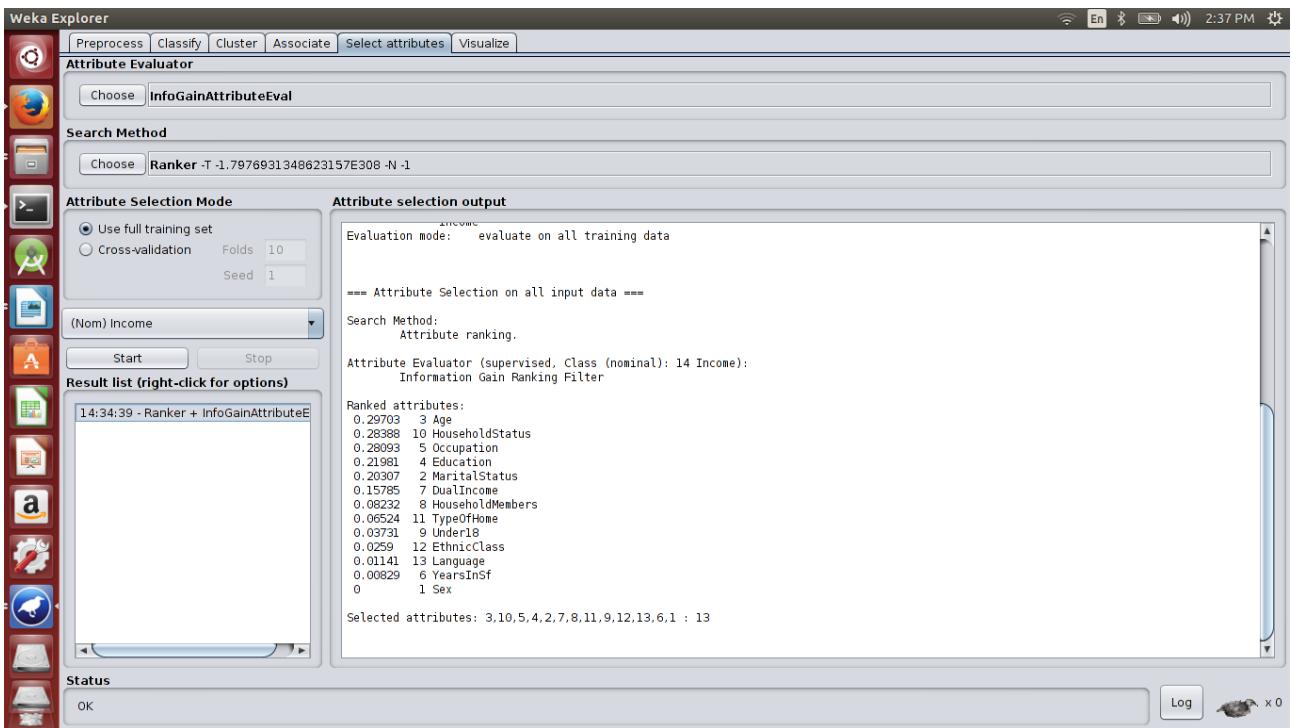
## 6)Attribute Selection Feature of Weka:

Use an attribute evaluator and a ranker to rank all the attributes of the dataset.Omit the attributes one at a time that have lower ranks to see the predictive accuracy of the classification algorithm.Weights put by the rankers algorithms are different than those by the classification algorithms.

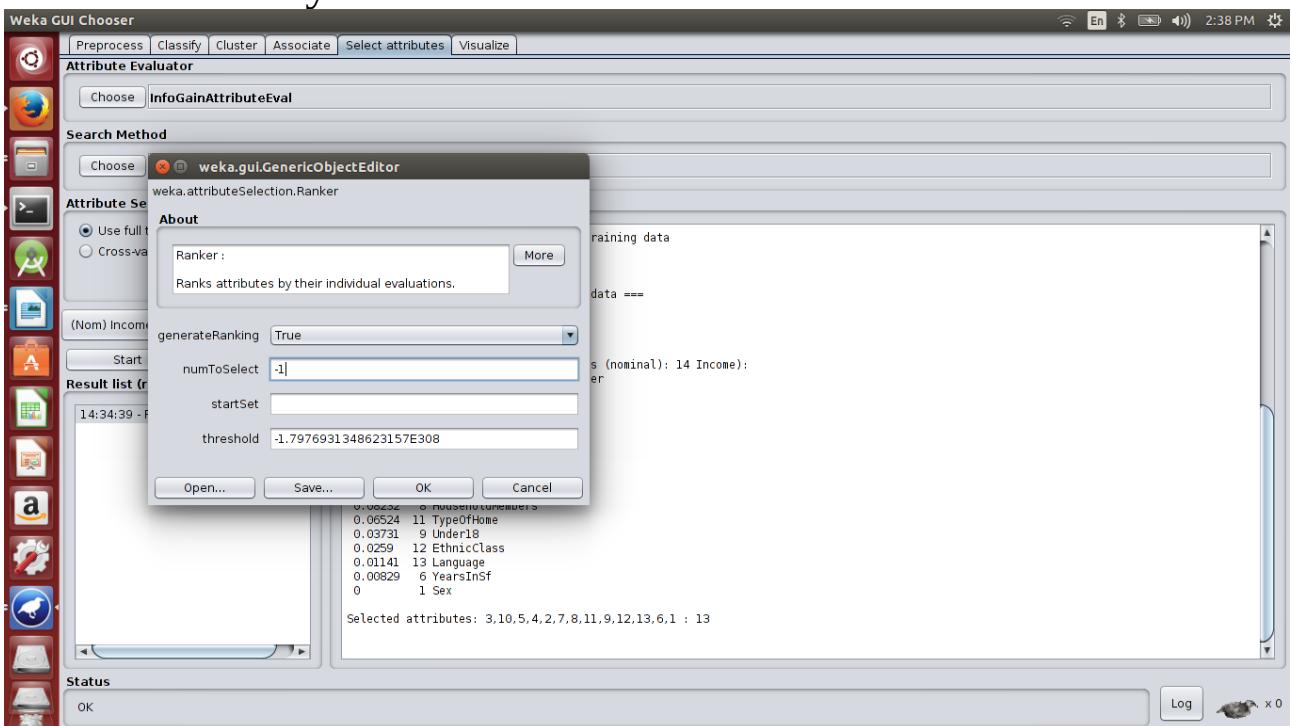
So, on selecting select attributes tab then selecting info gain attribute eval from attribute selection method gives the following alert that we need to use ranker's algo to use this (since default it is set to something else).



Click yes, and proceed further.clicking start button we obtain this:



So, the above image shows all the attributes ranks.  
And the number of attributes we want to select from attribute vector can always be defined.



Num to select if set to -1 it displays all the ranks. we can change it in case if we need only top few ranks.  
If all the 13 attributes are used correctly classified instances are 2828.

**Weka Explorer**

**Classifier**

Choose **NaiveBayes**

**Test options**

- Use training set
- Supplied test set Set...
- Cross-validation Folds 10
- Percentage split % 66
- More options...

(Nom) Income

Start Stop

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

- 14:48:45 - bayes.NaiveBayes
- 14:49:39 - bayes.NaiveBayes
- 14:53:29 - bayes.NaiveBayes
- 14:54:06 - bayes.NaiveBayes
- 14:55:14 - bayes.NaiveBayes
- 14:56:02 - bayes.NaiveBayes
- 14:56:26 - bayes.NaiveBayes
- 14:57:23 - bayes.NaiveBayes
- 14:58:02 - bayes.NaiveBayes
- 14:58:22 - bayes.NaiveBayes
- 14:58:35 - bayes.NaiveBayes
- 14:58:46 - bayes.NaiveBayes

**Classifier output**

```
==== Evaluation on training set ====
Time taken to test model on training data: 0.36 seconds
==== Summary ====
Correctly Classified Instances      2828      31.4467 %
Incorrectly Classified Instances   6165      68.5533 %
Kappa statistic                   0.2059
Mean absolute error               0.1632
Root mean squared error          0.3017
Relative absolute error           83.7339 %
Root relative squared error     96.6477 %
Total Number of Instances        8993

==== Detailed Accuracy By Class ====

```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
a	0.752	0.157	0.542	0.752	0.637	0.534	0.884	0.664	1
b	0.147	0.049	0.227	0.147	0.177	0.119	0.737	0.192	2
c	0.043	0.022	0.136	0.043	0.066	0.037	0.701	0.131	3
d	0.365	0.172	0.274	0.365	0.395	0.141	0.685	0.165	4
e	0.026	0.014	0.142	0.026	0.040	0.028	0.672	0.129	5
f	0.084	0.056	0.174	0.084	0.113	0.039	0.646	0.174	6
g	0.030	0.013	0.223	0.030	0.053	0.045	0.666	0.173	7
h	0.542	0.296	0.265	0.542	0.356	0.221	0.727	0.293	8
i	0.255	0.055	0.336	0.255	0.290	0.227	0.753	0.277	9
Weighted Avg.	0.314	0.103	0.279	0.314	0.267	0.196	0.734	0.290	

```
==== Confusion Matrix ====

```

	a	b	c	d	e	f	g	h	i	-- classified as
a	1313	161	12	143	9	46	9	44	8	a = 1
b	250	114	19	219	13	52	9	85	14	b = 2
c	...	...	...	...	...	...	...	...	...	

**Status**

OK Log x 0

when we remove few attributes by selecting them in preprocess as below:

**Weka Explorer**

**Preprocess**

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

**Filter**

Choose **ReplaceMissingValues**

**Current relation**

Relation: marketing-weka.filters.unsupervised.attribute.ReplaceMiss... Attributes: 14 Instances: 8993 Sum of weights: 8993

**Attributes**

All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> Sex
2	<input type="checkbox"/> MaritalStatus
3	<input type="checkbox"/> Age
4	<input type="checkbox"/> Education
5	<input type="checkbox"/> Occupation
6	<input checked="" type="checkbox"/> YearsInsf
7	<input type="checkbox"/> DualIncome
8	<input type="checkbox"/> HouseholdMembers
9	<input type="checkbox"/> Under18
10	<input type="checkbox"/> HouseholdStatus
11	<input type="checkbox"/> TypeOfHome
12	<input type="checkbox"/> EthnicClass
13	<input checked="" type="checkbox"/> Language
14	<input type="checkbox"/> Income

Remove

**Selected attribute**

Name: Language  
Missing: 0 (0%) Distinct: 4 Type: Numeric Unique: 0 (0%)

Statistic	Value
Minimum	1
Maximum	3
Mean	1.128
StdDev	0.406

**Class: Income (Nom)** Visualize All

**Status**

OK Log x 0

Now correctly classified increased to 2853

**Weka Explorer**

**Classifier**

Choose **NaiveBayes**

**Test options**

- Use training set
- Supplied test set Set...
- Cross-validation Folds 10
- Percentage split % 66
- More options...

(Nom) Income

Start Stop

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

- 14:48:45 - bayes.NaiveBayes
- 14:49:39 - bayes.NaiveBayes
- 14:53:29 - bayes.NaiveBayes
- 14:54:06 - bayes.NaiveBayes
- 14:55:14 - bayes.NaiveBayes
- 14:56:02 - bayes.NaiveBayes
- 14:56:26 - bayes.NaiveBayes
- 14:57:23 - bayes.NaiveBayes
- 14:58:02 - bayes.NaiveBayes
- 14:58:22 - bayes.NaiveBayes
- 14:58:35 - bayes.NaiveBayes
- 14:58:46 - bayes.NaiveBayes
- 14:59:47 - bayes.NaiveBayes

**Classifier output**

```
==== Evaluation on training set ====
Time taken to test model on training data: 0.3 seconds
==== Summary ====
Correctly Classified Instances      2853      31.7247 %
Incorrectly Classified Instances   6140      68.2753 %
Kappa statistic                   0.2088
Mean absolute error               0.1632
Root mean squared error          0.3007
Relative absolute error           83.7834 %
Root relative squared error     96.3195 %
Total Number of Instances        8993

==== Detailed Accuracy By Class ====

```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
a	0.768	0.157	0.541	0.768	0.635	0.541	0.884	0.668	1
b	0.124	0.036	0.245	0.124	0.165	0.121	0.742	0.196	2
c	0.019	0.007	0.178	0.019	0.035	0.026	0.709	0.141	3
d	0.365	0.172	0.274	0.365	0.395	0.140	0.685	0.165	4
e	0.004	0.001	0.214	0.004	0.008	0.019	0.672	0.127	5
f	0.084	0.058	0.162	0.084	0.107	0.039	0.653	0.175	6
g	0.137	0.026	0.250	0.137	0.271	0.137	0.731	0.296	7
h	0.514	0.244	0.264	0.514	0.349	0.211	0.731	0.294	8
i	0.295	0.066	0.329	0.295	0.311	0.241	0.755	0.278	9
Weighted Avg.	0.317	0.103	0.290	0.317	0.261	0.196	0.737	0.291	

```
==== Confusion Matrix ====

```

	a	b	c	d	e	f	g	h	i	-- classified as
a	1341	126	11	153	1	51	4	50	8	a = 1
b	271	96	14	254	2	63	2	79	14	b = 2
c	...	...	...	...	...	...	...	...	...	

**Status**

OK Log x 0

So,in this way by using classify,select attribute tab and preprocess tab and by using a particular classifying algorithm and particular feature selection algoritm we end up with only the useful attributes.Also we can use Wrapper method in which we can test with all possible subsets of attributes.