

Q1

1. What's the information gain for the 'Price' attribute? Please show your calculation.

Ans:

In total, P: 7, NP: 5

$$\text{Info}(D) = I(7, 5) = -7/12 \cdot \log(7/12) - 5/12 \cdot \log(5/12) = 0.97987$$

	Low	Medium	High
P	4	2	1
NP	1	1	3

$$\begin{aligned} \text{Info\_price}(D) &= 5/12 \cdot I(4, 1) + 3/12 \cdot I(2, 1) + 4/12 \cdot I(1, 3) \\ &= 5/12 \cdot (-4/5 \cdot \log_2(4/5) - 1/5 \cdot \log_2(1/5)) + 3/12 \cdot (-1/3 \cdot \log_2(1/3) - 2/3 \cdot \log_2(2/3)) \\ &\quad + 4/12 \cdot (-1/4 \cdot \log_2(1/4) - 3/4 \cdot \log_2(3/4)) \\ &= 0.800803 \end{aligned}$$

$$\text{So Gain(Price)} = \text{Info}(D) - \text{Info\_price}(D) = 0.179067$$

2. Now suppose we want to use Gini Index as attribute selection measure. What's the Gini index for the attribute Parking? What's the reduction in impurity in terms of Gini Index? Please show your calculation.

Ans: Parking partitions D into {No} and {Available}.

	No	Available
P	2	5
NP	2	3

$$\begin{aligned} \text{Gini}(D) &= 1 - (5/12)^2 - (7/12)^2 = 0.48611 \\ \text{Gini\_parking}(D) &= 4/12 \cdot \text{Gini}(2, 2) + 8/12 \cdot \text{Gini}(5, 3) \\ &= \frac{4}{12} \cdot \left(1 - \frac{1^2}{2} - \frac{1^2}{2}\right) + \frac{8}{12} \cdot \left(1 - \frac{3^2}{8} - \frac{5^2}{8}\right) = \frac{23}{48} \approx 0.47917 \end{aligned}$$

$$\text{Reduction in impurity} = \text{Gini}(D) - \text{Gini\_parking}(D) = 1/144 = 0.00694$$

3. Based on the training data, we want to construct a Naive Bayes classifier. Please estimate the following terms (No smoothing is required, and please show your calculation):

Ans:

$$\text{a) } \Pr(\text{Popularity} = \text{'P'}) = 7/12; \text{ and } \Pr(\text{Popularity} = \text{'N'}) = 5/12$$

$$\begin{aligned} \text{b) } \Pr(\text{Price} = \text{'Low'}, \text{Parking} = \text{'Available'}, \text{Cuisine} = \text{'Mexican'} \mid \text{Popularity} = \text{'P'}) \\ = \Pr(\text{Low} \mid \text{P}) \cdot \Pr(\text{Available} \mid \text{P}) \cdot \Pr(\text{Mexican} \mid \text{P}) = 4/7 \cdot 5/7 \cdot 2/7 = 40/343 \end{aligned}$$

$$\begin{aligned} \text{c) } \Pr(\text{Price} = \text{'Low'}, \text{Parking} = \text{'Available'}, \text{Cuisine} = \text{'Mexican'} \mid \text{Popularity} = \text{'N'}) \\ = \Pr(\text{Low} \mid \text{NP}) \cdot \Pr(\text{Available} \mid \text{NP}) \cdot \Pr(\text{Mexican} \mid \text{NP}) = 1/5 \cdot 3/5 \cdot 2/5 = 6/125 \end{aligned}$$

4. Suppose a restaurant has the values: Price = 'Low', Parking = 'Available', Cuisine = 'Mexican'. Based on the calculation in 3, is this restaurant classified as popular?

$$P1 = P(P \mid \text{Price} = \text{'Low'}, \text{Parking} = \text{'Available'}, \text{Cuisine} = \text{'Mexican'})$$

$$= \frac{P(\text{Price} = \text{low}, \text{Parking} = \text{'Available'}, \text{Cuisine} = \text{'Mexican'} \mid P) * P(P)}{P(\text{Price} = \text{'Low'}, \text{Parking} = \text{'Available'}, \text{Cuisine} = \text{'Mexican'})}$$

$$= \frac{\left(\frac{40}{343} * \frac{7}{12}\right)}{P(\text{Price} = \text{'Low'}, \text{Parking} = \text{'Available'}, \text{Cuisine} = \text{'Mexican'})}$$

$$P2 = P(NP \mid \text{Price} = \text{'Low'}, \text{Parking} = \text{'Available'}, \text{Cuisine} = \text{'Mexican'})$$

$$= \frac{P(\text{Price} = \text{low}, \text{Parking} = \text{'Available'}, \text{Cuisine} = \text{'Mexican'} \mid NP) * P(NP)}{P(\text{Price} = \text{'Low'}, \text{Parking} = \text{'Available'}, \text{Cuisine} = \text{'Mexican'})}$$

$$= \frac{\left(\frac{6}{125} * \frac{5}{12}\right)}{P(\text{Price} = \text{'Low'}, \text{Parking} = \text{'Available'}, \text{Cuisine} = \text{'Mexican'})}$$

As  $P1 > P2$ , it is classified as Popular.

Q2:

1. According to training data, we know that red is +1, blue is -1.

For each test data, we compute and compare the distance to all training points to get the K=3 nearest neighbors and assign labels according to it.

train				Euclidean dist <sup>2</sup> (Test to train)			
x1	x2	y		1	2	3	4
1	0.5	1		7.73	2.5	4.25	0.29
2	1.2	1		2.74	0.29	1.94	0.68
2.5	2	1		0.53	1	1.25	2.69
3	2	1		0.58	1.25	2.5	4.24
1.5	2	-1		1.93	2	0.25	1.09
2.3	3	-1		0.25	4.04	0.89	5.21
1.2	1.9	-1		2.89	2.5	0.45	0.81
0.8	1	-1		6.5	2.89	2.74	0.16
			Label(y)	1	1	-1	1

Those 3 yellow cells in each column corresponds to the 3 nearest points to current test point. We get the classification labels shown in green. The error is the 4<sup>th</sup> point label, which should be -1 in this case, with error rate = 1/4

2.  $\text{sign}(w^T x) = \text{sign}(w_0 + w_1 x_1 + w_2 x_2) = \text{sign}(1 + 0.5 * 0.8 - 1 * 1) = \text{sign}(0.4) = 1$  which is not -1. So it is NOT correctly classified.

Adjustment: Update  $w = w + \eta xy$

$$w_0 = w_0 - 0.1 * 1 = 0.9$$

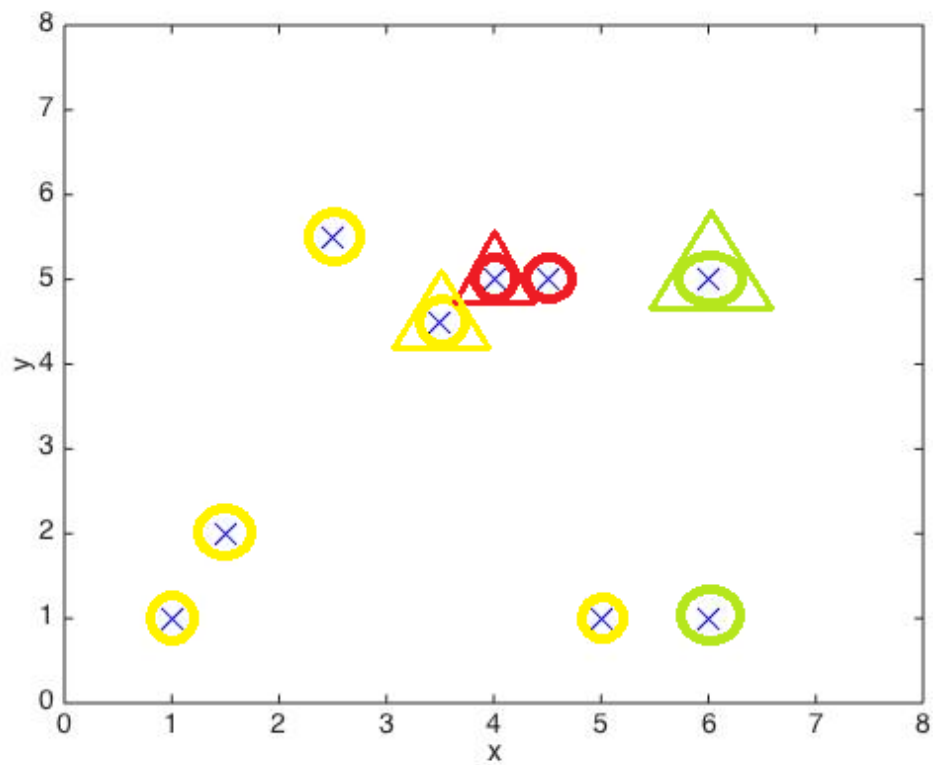
$$w_1 = w_1 - 0.1 * 0.8 = 0.42$$

$$w_2 = w_2 - 0.1 * 1 = -1.1$$

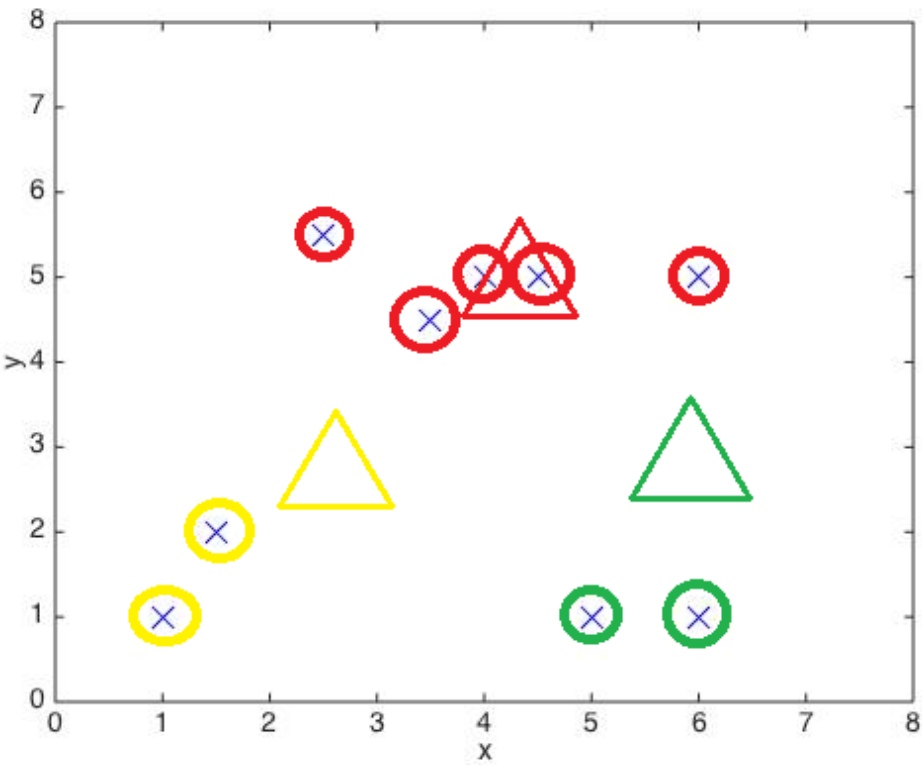
Q3:

1. Note: Triangles represent cluster centers

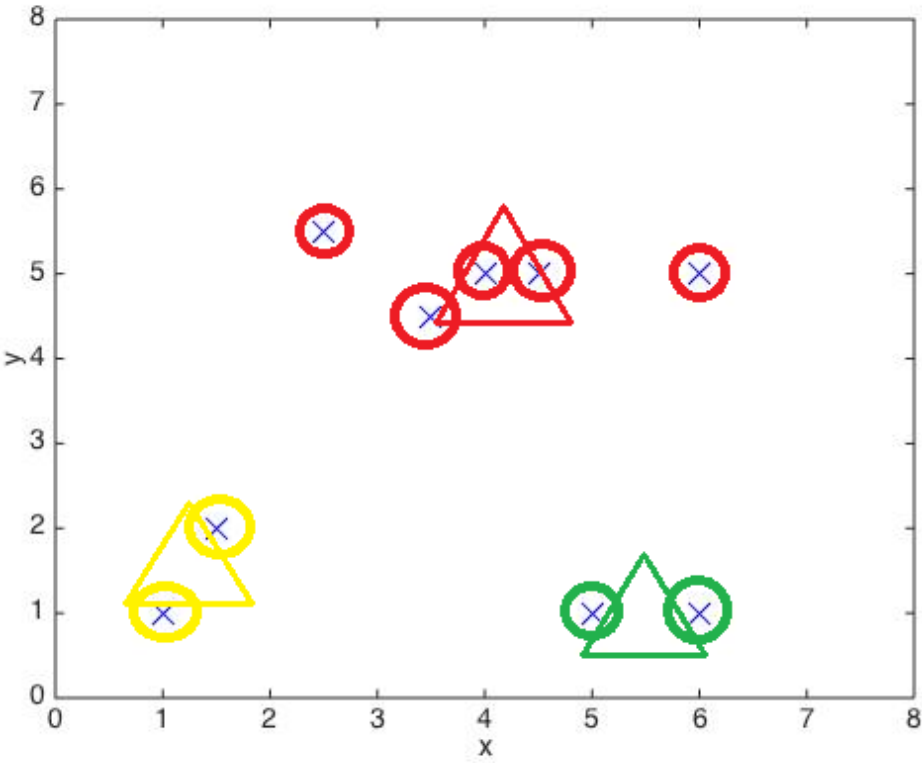
Initial	Center Cord x	Center Cord y
Cluster 1	4	5
Cluster 2	3.5	4.5
Cluster 3	6	5



Round 2:	Center Cord x	Center Cord y
Cluster 1	4.25	5
Cluster 2	2.7	2.8
Cluster 3	6	3

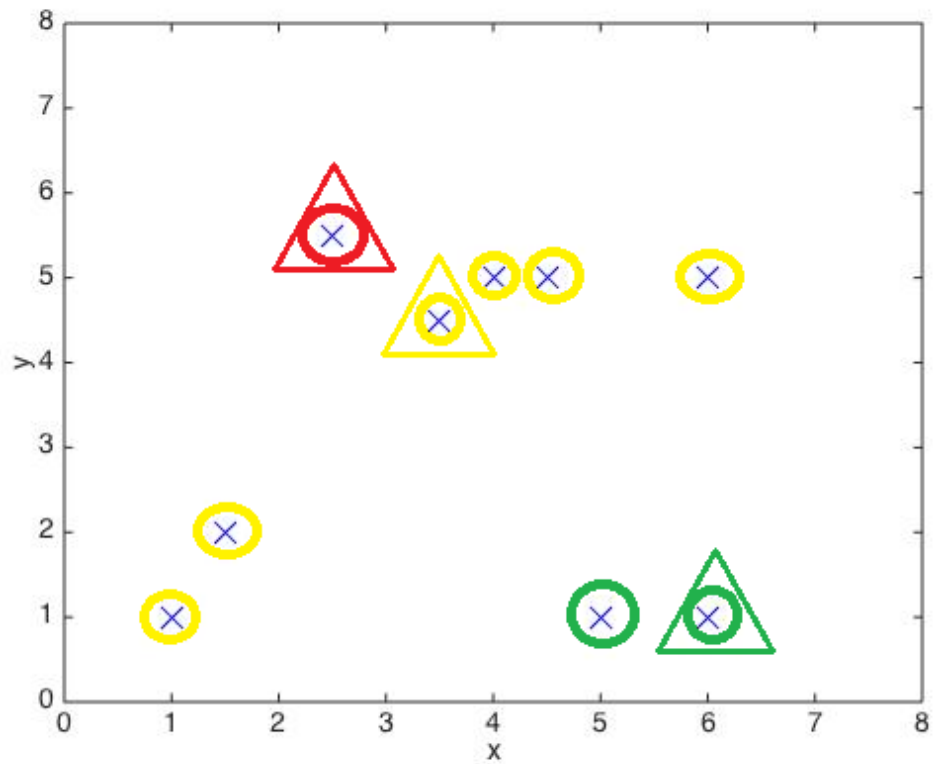


Round 3:	Center Cord x	Center Cord y
Cluster 1	4.1	5
Cluster 2	1.25	1.5
Cluster 3	5.5	1

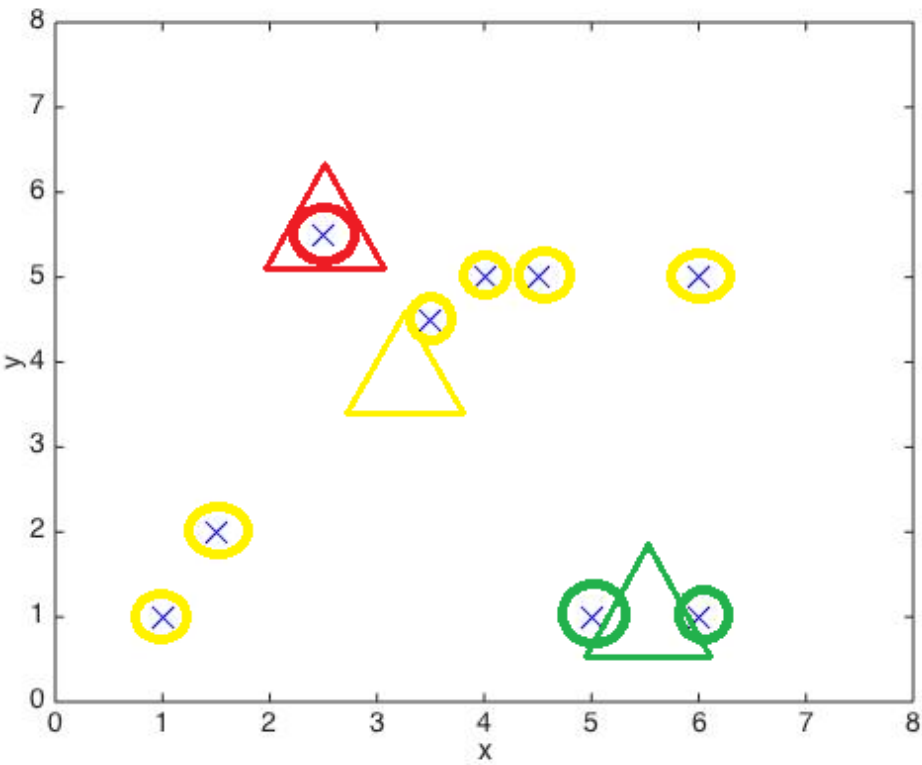


2. Note: Triangles represent cluster centers

Initial	Center Cord x	Center Cord y
Cluster 1	2.5	5.5
Cluster 2	3.5	4.5
Cluster 3	6	1



Round 2	Center Cord x	Center Cord y
Cluster 1	2.5	5.5
Cluster 2	3.416667	3.75
Cluster 3	5.5	1





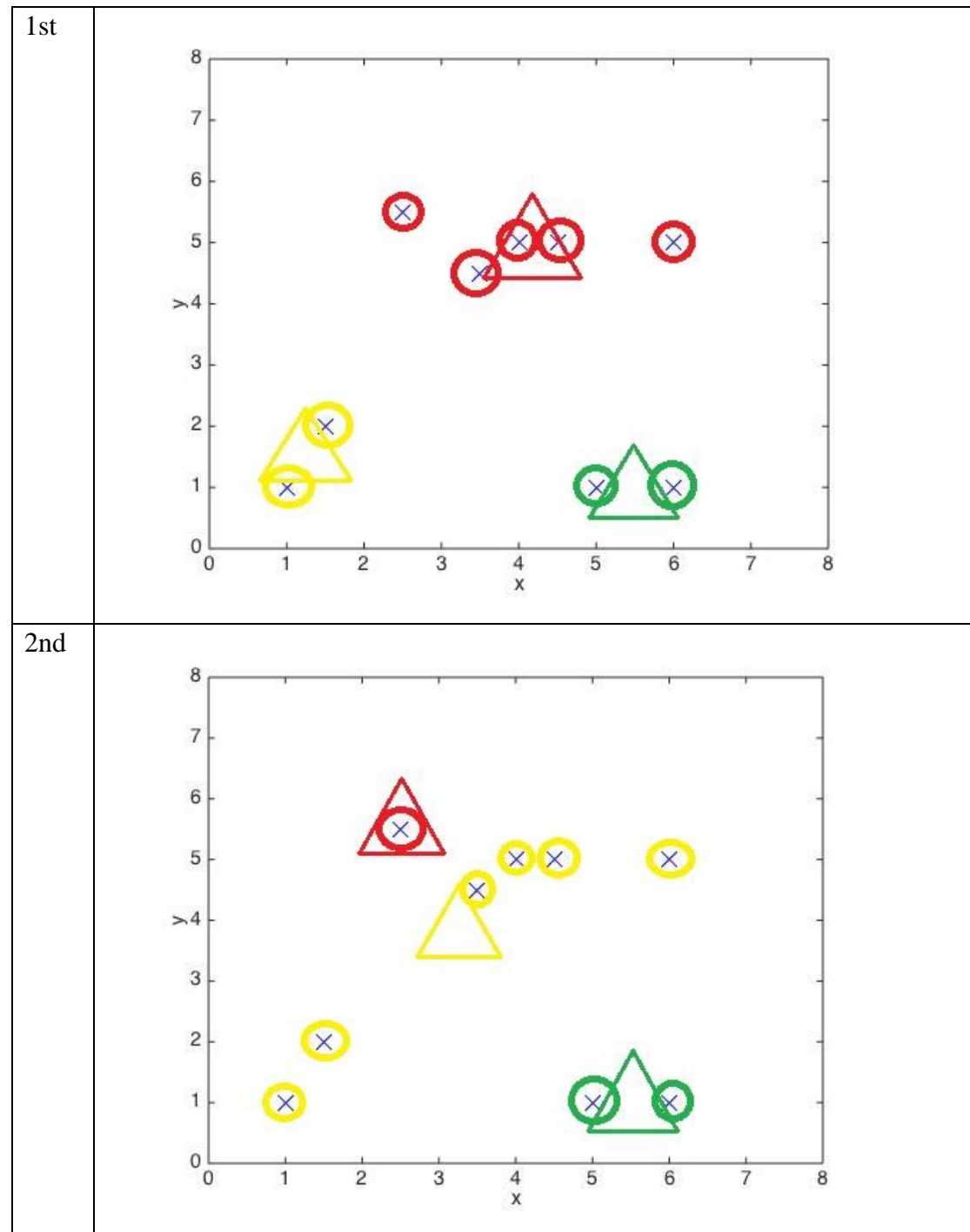
3.

1) They are different because:

It terminates at local optimal. As for the 2<sup>nd</sup>, the cluster graphs is more distorted by the outliers.

2) The first one is better. As can be seen from the two graphs below:

The 1<sup>st</sup> one produced high quality clusters with more cohesive within clusters and distinctive between clusters than the 2<sup>nd</sup> one.

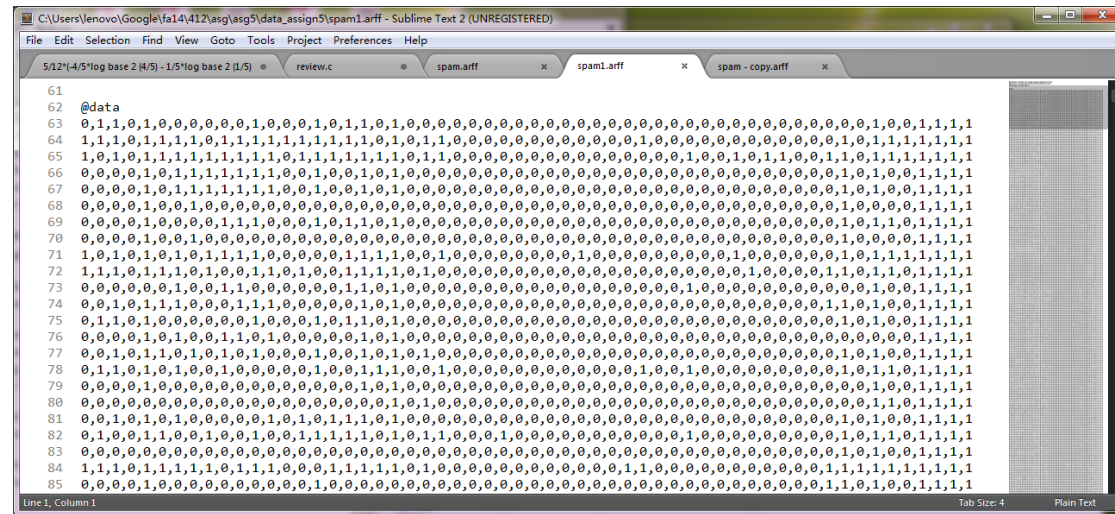


3) Quality measurement criteria: We can calculate the total distance among each clusters and add them up. The smaller this value is, the better the cluster is. For this example, the 1<sup>st</sup> case is better, because the 2<sup>nd</sup> one have spread yellow cluster.

4) We can try to exclude outliers from the initial cluster center or try out multiple starting points and choosing the clustering with lowest cost (just like the above process).

## Mini-Mp:

### 1. Binary Attribute:



### 2. J48 Result:

```
=== Evaluation on test split ===  
=== Summary ===
```

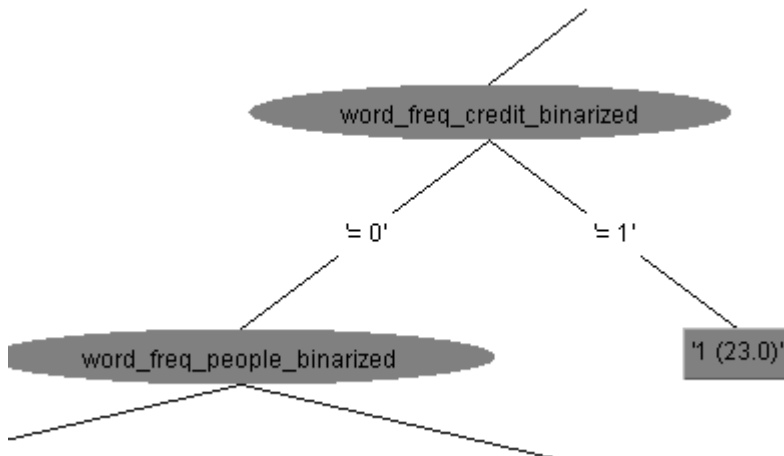
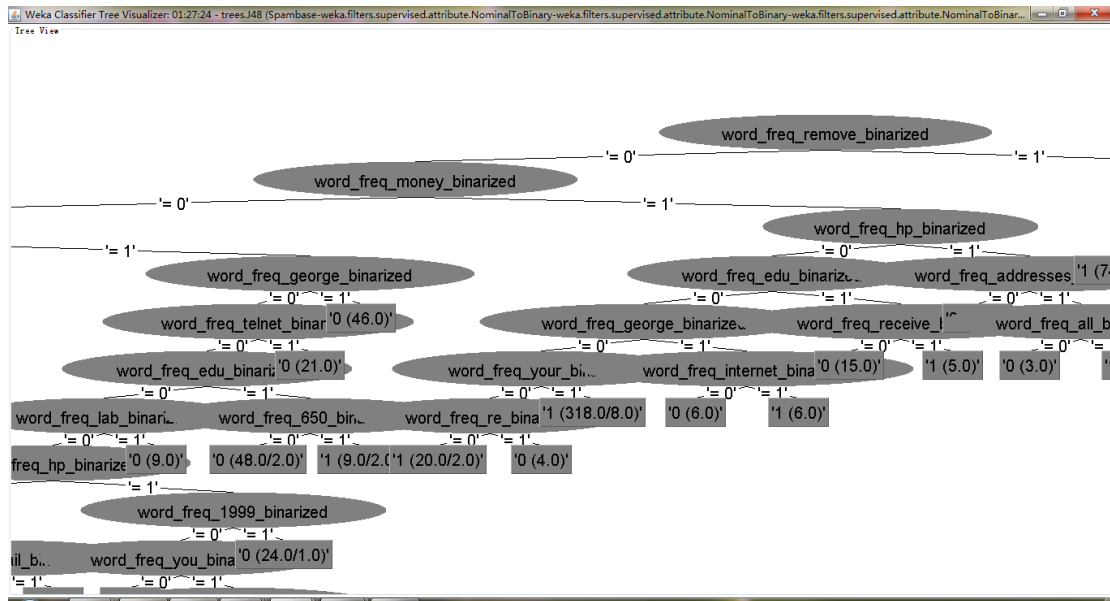
Correctly Classified Instances	1426	91.1765 %
Incorrectly Classified Instances	138	8.8235 %
Kappa statistic	0.8131	
Mean absolute error	0.1229	
Root mean squared error	0.2782	
Relative absolute error	25.7356 %	
Root relative squared error	56.9163 %	
Total Number of Instances	1564	

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

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
	0.949	0.146	0.909	0.949	0.929	0.937	0
	0.854	0.051	0.917	0.854	0.884	0.937	1
Weighted Avg.	0.912	0.108	0.912	0.912	0.911	0.937	

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

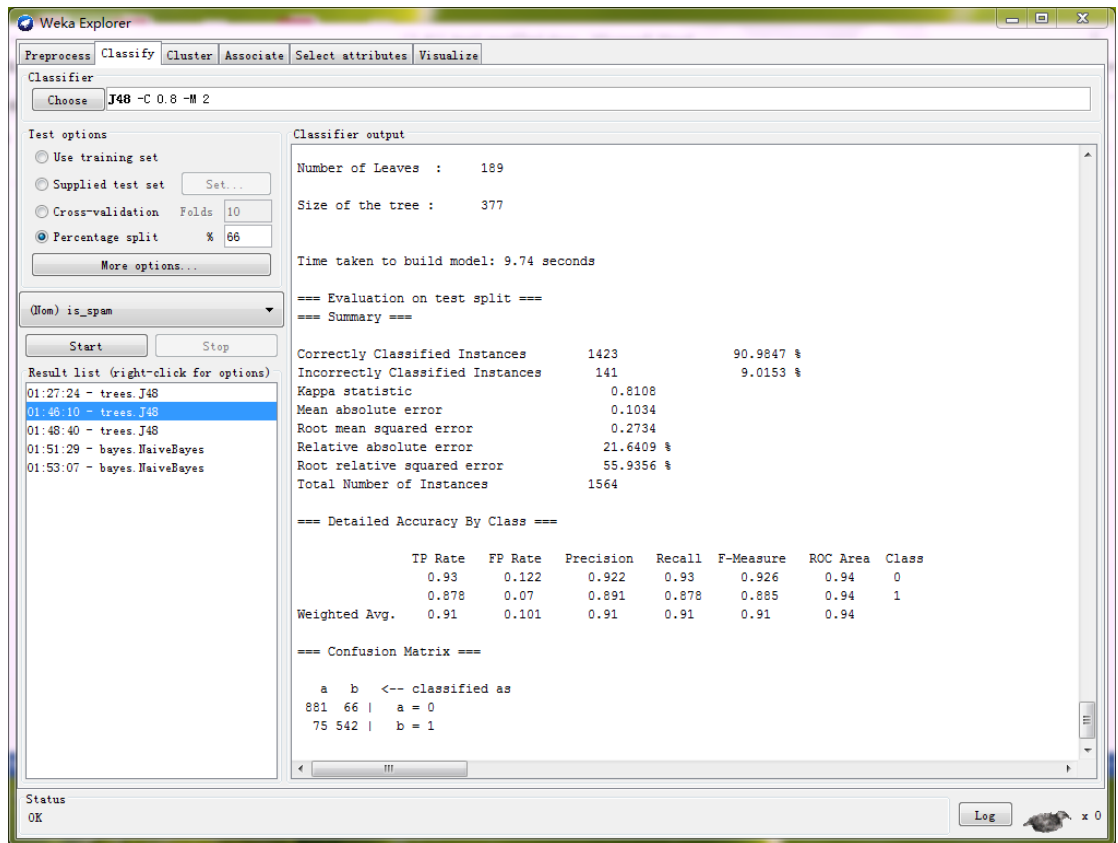
```
  a   b  <-- classified as  
899  48 |   a = 0  
 90 527 |   b = 1
```



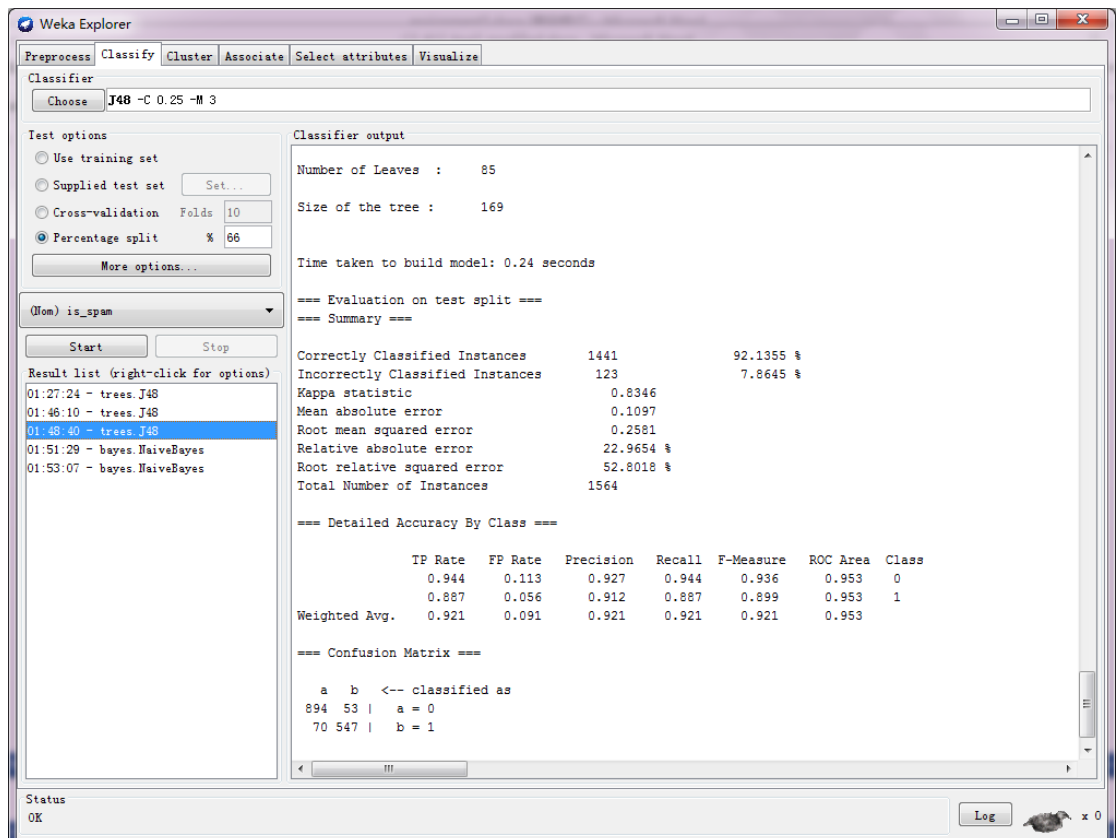
Interesting Rule: Word freq with credit turns out more likely to be classified as spam

### 3. Play with weka:

- 1) **Confidencefactor** is used for pruning. I chose 0.8, larger than default 0.25, then, the tree get much bigger.



- 2) **MinnumberObj** is to ensure the instances of leaves is bigger than this number. After I change this from 2 to 3, the leaves decreases with bigger instance.



#### 4. Naïve Bayes:

The screenshot shows the Weka Explorer interface with the 'Classify' tab selected. The 'Classifier' dropdown is set to 'NaiveBayes'. Under 'Test options', 'Percentage split' is selected with a value of 86%. The 'Result list' on the left shows several entries, with '01:53:07 - bayes.NaiveBayes' selected. The 'Classifier output' pane displays the following results:

=== Evaluation on test split ===  
=== Summary ===

Correctly Classified Instances	1380	88.2353 %
Incorrectly Classified Instances	184	11.7647 %
Kappa statistic	0.7508	
Mean absolute error	0.1202	
Root mean squared error	0.3174	
Relative absolute error	25.1712 %	
Root relative squared error	64.9514 %	
Total Number of Instances	1564	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
	0.925	0.183	0.886	0.925	0.905	0.952	0
	0.817	0.075	0.877	0.817	0.846	0.952	1
Weighted Avg.	0.882	0.14	0.882	0.882	0.882	0.952	

=== Confusion Matrix ===

```
a  b  <-- classified as
876  71 |  a = 0
113 504 |  b = 1
```

The status bar at the bottom shows 'Status OK' and a 'Log' button.