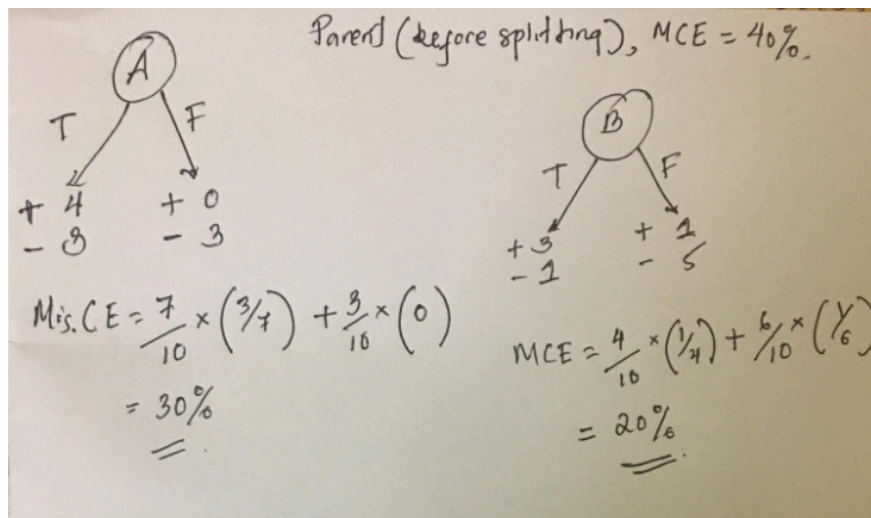


Data Mining Assignment 4

- 1) Read Chapter 4 (all sections) and Chapter 5 (Sections 5.2, 5.5, 5.6 and 5.7).
- 2) Consider the following data set for a binary class problem.

A	B	Class Label
T	F	+
T	T	+
T	T	+
T	F	-
T	T	+
F	F	-
F	F	-
F	F	-
T	T	-
T	F	-

Calculate the misclassification error rate when splitting on A and B to determine the best split. Which of these splits considered is the best according to misclassification error rate?



According to the misclassification error rate, splitting on attribute B is considered to be the best split over attribute A.

3) Consider the training examples shown below for a binary classification problem.

Instance	a_1	a_2	a_3	Target Class
1	T	T	1.0	+
2	T	T	6.0	+
3	T	F	5.0	-
4	F	F	4.0	+
5	F	T	7.0	-
6	F	T	3.0	-
7	F	F	8.0	-
8	T	F	7.0	+
9	F	T	5.0	-

For a_3 , which is a continuous attribute compute misclassification error rate for every possible split to determine the best split. Which of these splits considered is the best according to misclassification error rate?

a_3	class	split position	MCE
1.0	+	2.0	0.33
3.0	-	3.5	0.44
4.0	+	4.5	0.33
5.0	-	5.5	0.44
6.0	+	6.5	0.44
7.0	-	7.5	0.44
8.0	-		

The best split for a_3 occurs at split point equals to 2 and 4.5.

4) The file http://www-stat.wharton.upenn.edu/~dmease/rpart_text_example.txt gives an example of text output for a tree fit using the rpart() function in R from the library rpart. Use this tree to predict the class labels for the 10 observations in the test data http://www-stat.wharton.upenn.edu/~dmease/test_data.csv linked here. Do this manually - do not use R or any software.

Age	Number	Start		Predicted
middle	5	10	absent	Present
young	2	17	absent	Absent
old	10	6	present	Present
young	2	17	absent	Absent
old	4	15	absent	absent
middle	5	15	absent	absent
young	3	13	absent	absent
old	5	8	present	Present
young	7	9	absent	absent
middle	3	13	absent	absent

5) I split the popular sonar data set into a training set (http://www-stat.wharton.upenn.edu/~dmease/sonar_train.csv) and a test set (http://www-stat.wharton.upenn.edu/~dmease/sonar_test.csv).

Use R to compute the misclassification error rate on the test set

when training on the training set for a tree of depth 5 using all the default values except control=rpart.control(minsplit=0,minbucket=0,cp=-1, maxcompete=0, maxsurrogate=0, usesurrogate=0, xval=0,maxdepth=5).

Remember that the 61st column is the response and the other 60 columns are the predictors.

```

> setwd("/Users/pemayangdon/Desktop/DS/Specialization/2. DataMining/data Mining
Assignments/DM Assignment4")
> sonar_train <- read.csv("sonar_train.csv", header=F)
> sonar_test <- read.csv("sonar_test.csv", header=T)
> y_train <- as.factor(sonar_train[,61])
> x_train <- sonar_train[,1:60]
> library(rpart)
> model <- rpart(y_train~.,x_train,control=rpart.control(minsplit=0,
minbucket=0,cp=-1,maxcompete=0,maxsurrogate=0,usesurrogate=0,xval=0,maxdepth=5))
> 1 - sum(y_train==predict(model, x_train, type="class")) / length(y_train)
[1] 0.01538462

> 1 - sum(y_test == predict(model, x_test, type="class")) / length(y_test)
[1] 0.2564103

```

Hence, the misclassification error rate on test data ~ 26%

6) Do Chapter 5 textbook problem #17 (parts a and c only) on pages 322-323. Note that there is a typo in part c - it should read "Repeat the analysis for part (b)". We will do part b in class.

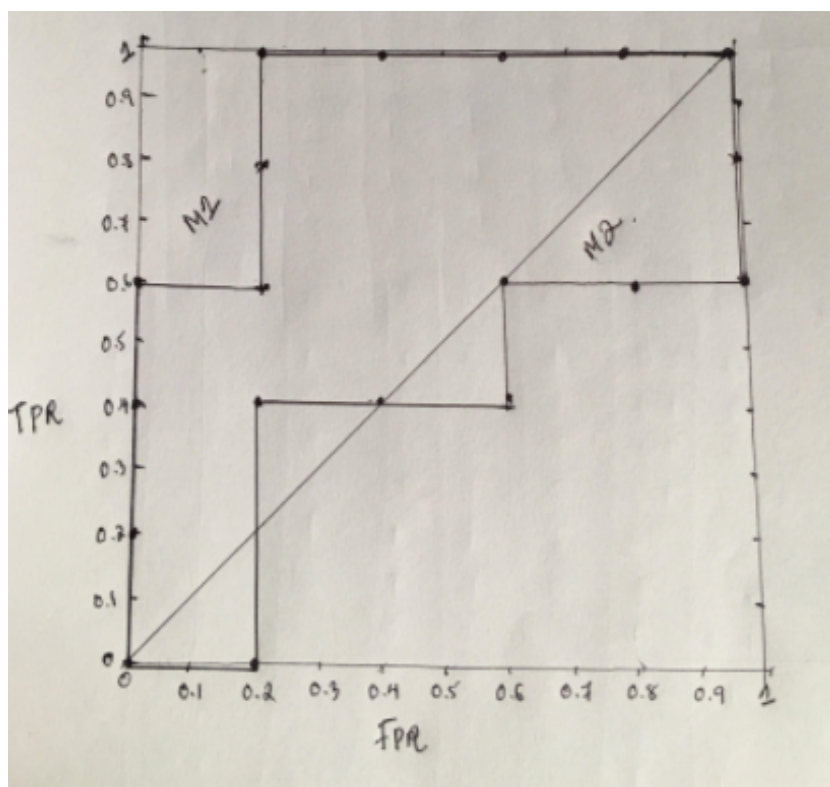
Q. You are asked to evaluate the performance of two classification models, M1 and M2. The test set you have chosen contains 26 binary attributes, labeled as A through Z.

Table 5.14 shows the posterior probabilities obtained by applying the models to the test set. (Only the posterior probabilities for the positive class are shown). As this is a two-class problem, $P(-) = 1 - p(+)$ and $p(-|A, \dots, Z) = 1 - P(+|A, \dots, Z)$. Assume that we are mostly interested in detecting instances from the positive class.

Table 5.14. Posterior probabilities for Exercise 17.

Instance	True Class	$P(+ A, \dots, Z, M_1)$	$P(+ A, \dots, Z, M_2)$
1	+	0.73	0.61
2	+	0.69	0.03
3	-	0.44	0.68
4	-	0.55	0.31
5	+	0.67	0.45
6	+	0.47	0.09
7	-	0.08	0.38
8	-	0.15	0.05
9	+	0.45	0.01
10	-	0.35	0.04

(a) Plot the ROC curve for both M1 and M2. (You should plot them on the same graph) Which model do you think is better? Explain your reasons.



M1 is far better than the M2, as its area under ROC curve (AUC) is much larger than the area under ROC curve for M2.

(b) For model M1, suppose you choose the cutoff threshold to be $t = 0.5$. In other words, any test instances whose posterior probability is greater than t will be classified as a positive example. Compute the precision, recall, and F-measure for the model at this threshold value.

threshold = 0.5

		+	-
Actual class	+	3	2
	-	2	4

$$\text{Precision} = \frac{TP}{TP+FP} = \frac{3}{3+1} = 75\%$$

$$\text{Recall} = \frac{TP}{TP+FN} = \frac{3}{3+2} = 60\%$$

$$\text{F-measure} = \frac{2rp}{r+p} = \frac{2 \times 0.75 \times 0.6}{0.75 + 0.6} = 0.667$$

(c) Repeat the analysis for part (b) using the same cutoff threshold on model M2. Compare the F-measure results for both models. Which model is better? Are the results consistent with what you expect from the ROC curve?

threshold = 0.5

		+	-
Achal	+	1	4
class	-	2	4

Precision = $\frac{1}{2} = 50\%$
Recall = $\frac{1}{5} = 20\%$
F-measure = $\frac{2 \times 0.5 \times 0.2}{0.5 + 0.2} = 0.286$

M1 is better, as its F-measure is greater than the F-measure of M2.

Yes, the result is consistent with the ROC plot.

7) Compute the misclassification error on the training data for the Random Forest classifier to the last column of the sonar training data. Show your R code for doing this.

```
> library(randomForest)
> fit <- randomForest(x,y)
> 1 - sum(y==predict(fit,x)) / length(y)
[1] 0
> 1 - sum(y_test==predict(fit,x_test)) / length(y_test)
[1] 0.1282051
```

8) This question deals with sonar data

a) Use `knn()` for the k-nearest neighbor classifier for $k=5$ and $k=6$ to the last column of the sonar training data. Compute the misclassification error on the training data and also on the test data.

```
<
> library(class)
> train <- read.csv("sonar_train.csv", header=F)
> y <- as.factor(train[,61])
> x <- train[,1:60]
> model <- knn(x,x,y,k=5)
> 1 - sum(y==model)/length(y)
[1] 0.2076923
> test <- read.csv("sonar_test.csv", header=F)
> y_test <- as.factor(test[,61])
> x_test <- test[,1:60]
> model_test <- knn(x, x_test, y, k=5)
> 1 - sum(y_test == model_test) / length(y_test)
[1] 0.2307692
>
> model1 <- knn(x,x,y,k=6)
> 1 - sum(y==model1) / length(y)
[1] 0.2615385
> model1_test <- knn(x,x_test,y,k=6)
> 1 - sum(y_test == model1_test) / length(y_test)
[1] 0.2820513
```

b) Repeat part a using the exact same R code a few times. Explain why both the training errors and the test errors often change for $k=6$ but not for $k=5$. Hint: Read the help on the `knn` function if you do not know.


```
> model1_test <- knn(x,x_test,y,k=6)
> 1 - sum(y_test == model1_test) / length(y_test)
[1] 0.2948718
> model1 <- knn(x,x,y,k=6)
> 1 - sum(y==model1) / length(y)
[1] 0.1923077
> model_test <- knn(x, x_test, y, k=5)
> 1 - sum(y_test == model_test) / length(y_test)
[1] 0.2307692
> model <- knn(x,x,y,k=5)
> 1 - sum(y==model) / length(y)
[1] 0.2076923
```

Because the classification is decided by majority vote, with ties broken at random.