

**ST4061 - Statistical Methods for Machine Learning II**  
**ST6041 - Machine Learning and Statistical Analytics II**

## CA1 Answer document

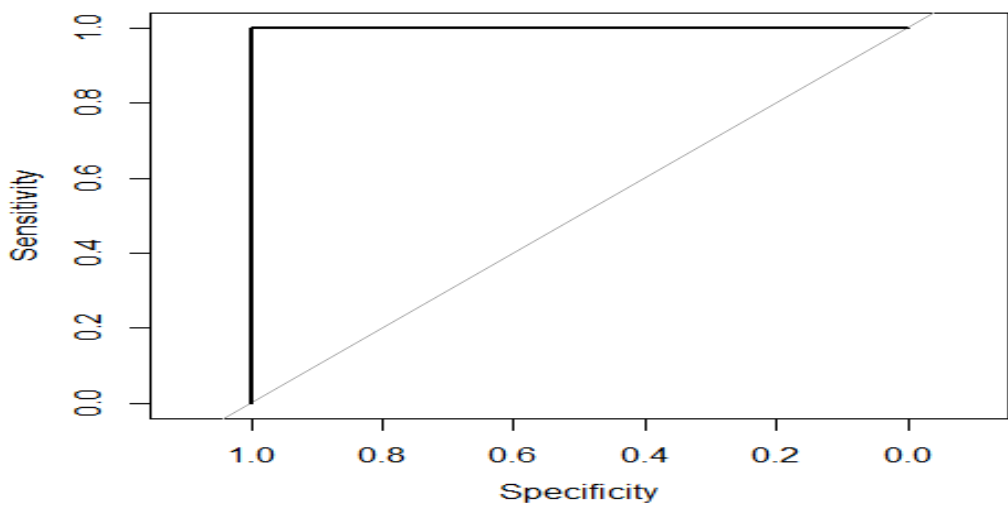
### Question 1

Question item	Answer					
(a)		<b>Length</b>	<b>Width</b>	<b>Leaf</b>	<b>Curve</b>	<b>Prediction</b>
	<b>Obs1</b>	4.5	2.3	1.3	0.3	<b>Young</b>
	<b>Obs2</b>	5.0	3.5	4.3	0.3	<b>Intermediate</b>
	<b>Obs3</b>	6.1	3.0	4.9	1.8	<b>Mature</b>
	<b>Obs4</b>	7.2	3.0	5.8	1.9	<b>Mature</b>
	<b>Obs5</b>	5.1	3.8	2.5	0.4	<b>Intermediate</b>
(b)	Misclassification rate = $100 \times (3/5) = 60\%$					
(c)	<p>i) Sum of all elements in the confusion matrix provided = <b>42</b></p> <p>ii) <math>100 \times (\text{Sum of all leading diagonal elements} / \text{Total}) = 100 \times ((8 + 12 + 9) / 42) = \mathbf{69.04\%}</math></p> <p>iii) <math>100 \times (8 / (8 + 3 + 1)) = \mathbf{66.67\%}</math></p> <p>iv) <math>100 \times ((0 + 3) / (0 + 3 + 9)) = \mathbf{25\%}</math></p>					

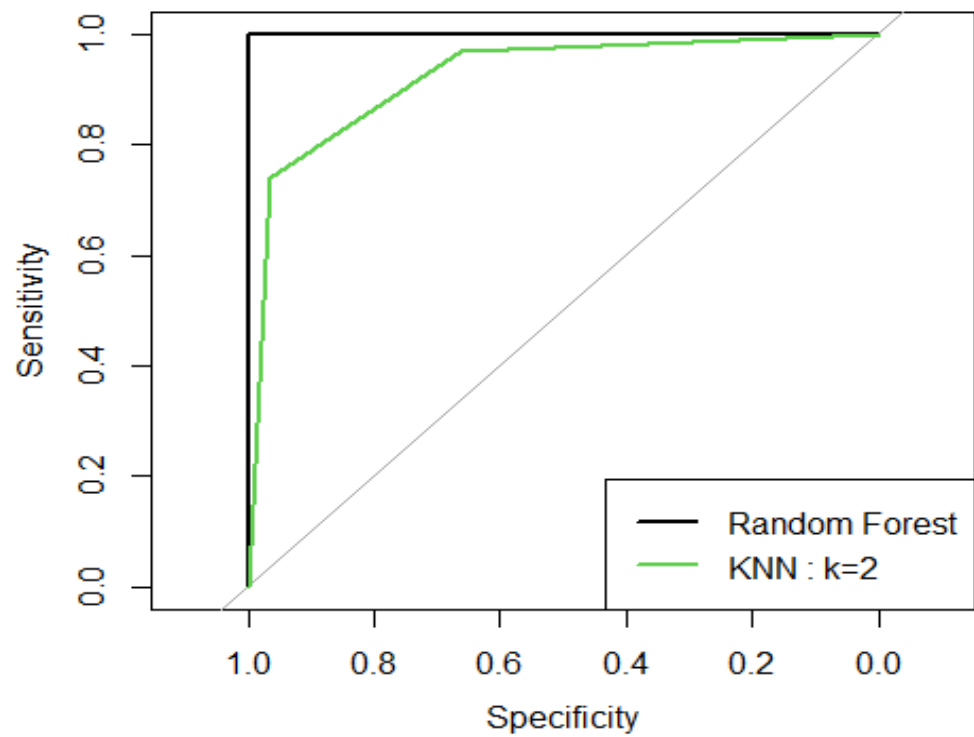
## Question 2

Question item	Answer
(a)	It's a <b>regression</b> problem as the dependent variable is salary (in log scale here) which is a range of numbers (continuous variable) and not a finite discrete set of values.
(b)	It is a <b>generalized linear model fit with elastic-net regularisation</b> . As $\alpha = 0.5$ , it is an exact 50-50 mix of L1 (Lasso) and L2 (Ridge) regression penalty applied to the loss (mean square error) function during the training of the model.
(c)	Type of Cross-Validation applied to mod2 is K-fold Cross-Validation and $K=10$ where dataset is divided into 10 folds using <code>folds=cut(1:n,K,labels=FALSE)</code> command. So, it's 10-fold CV. Also, <code>cv.glmnet()</code> command uses a 10-fold CV internally to compute the optimum lambda value and is passed as a parameter for each of the 'K' mod2 fits.

## Question 3

Question item	Answer
(a)	<p>Misclassification error rate from training data = <b>0%</b></p> <pre>&gt; prediction.conf</pre> <pre>rf.tree.preds Down Up Down 486 0 Up 0 514</pre>
(b)	 <p>Area under the curve: <b>1</b></p>

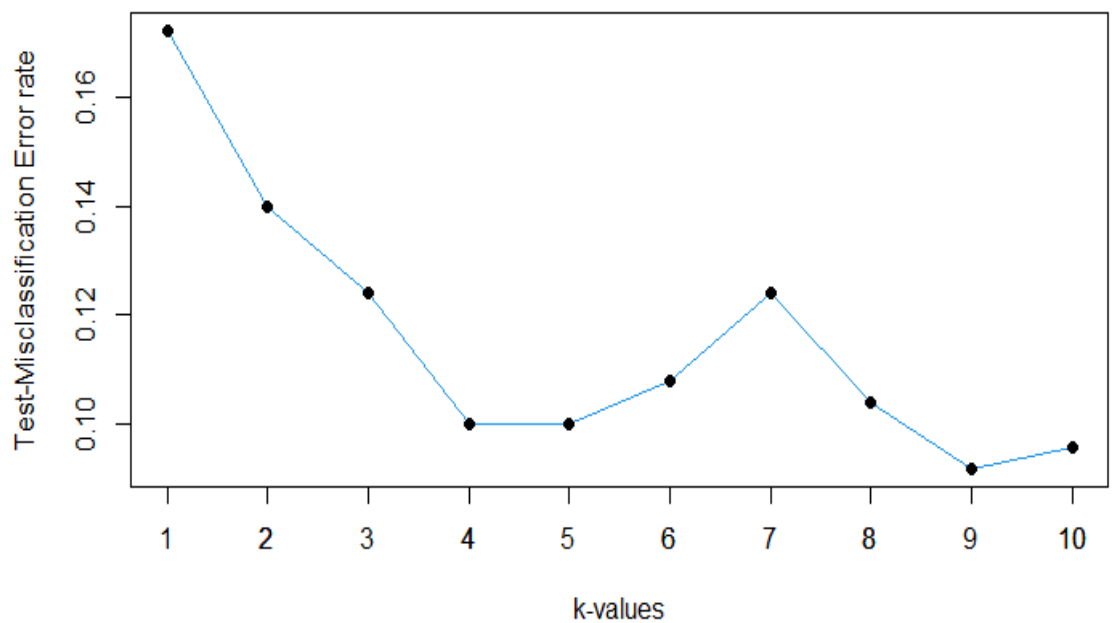
(c)



Area under the curve for KNN,  $k = 2$ : 0.9244

(d)

**Test-set Misclassification errors for KNN with k-values = 1:10**



### R code for Question 3

```
rm(list = ls())
require(ISLR)
require(class)
require(pROC)
library(randomForest)

x = Smarket[,-9]
y = Smarket$Direction
set.seed(4061)
train = sample(1:nrow(Smarket),1000)

## Question 3 i)

rf.tree = randomForest(y[train] ~ ., data = x[train,])
rf.tree
#summary(rf.tree)
rf.tree.preds = predict(rf.tree, x[train,], type = 'class')
prediction.conf = table(rf.tree.preds, y[train])
prediction.conf
missclass_rate = (1 - (sum(diag(prediction.conf))/sum(prediction.conf)))
missclass_rate

## Question 3 ii)
y_test_true = y[-train]
test_preds = predict(rf.tree, newdata=x[-train,], type='class')
#(rf.test.confusion = table(test_preds, y_test_true))
rftree.probs = predict(rf.tree, x[-train,], type="prob")
roc = roc(response=y_test_true, predictor=rftree.probs[,2])
auc = roc$auc
auc
plot(roc, col=1)

## Question 3 iii)
k = 2
knn.o = knn(x[train,], x[-train,], y[train], k)
knn.preds = as.numeric(knn.o == 'Up')
knn.p = attributes(knn(x[train,], x[-train,], y[train], k, prob=TRUE))$prob
new.probs = 1 - knn.p
final.knn.preds = ifelse(knn.preds == 1,knn.p, new.probs)
roc_knn = roc(y_test_true, final.knn.preds)
plot(roc_knn, add = TRUE, col = 75)
legend("bottomright", legend = c("Random Forest", "KNN : k=2"), col = c(1, 75), lty = 1,
lwd = 2)
auc_knn = roc_knn$auc
auc_knn
```

```

## Question 3 iv)
set.seed(4061)
M = 1000
train = sample(1:nrow(Smarket), M)

K = 10
test_class_errors = numeric(K)*NA
for(k in 1:K) {
  knn.o = knn(x[train,], x[-train,], y[train], k)
  confusion_mat = table(knn.o, y[-train])
  test_class_errors[k] = (1 - (sum(diag(confusion_mat))/sum(confusion_mat)))
}
test_class_errors
plot(seq(1:K), test_class_errors, xlim = c(1,10),
     xlab = "k-values", ylab = "Test-Misclassification Error rate",
     main = paste("Test-set Misclassification errors for KNN with k-values = 1:",K,sep=""),
     col = 4, type = 'l')
points(seq(1:K), test_class_errors, col=1, pch=20, cex = 1.4)
axis(side = 1, at = seq(1, 10, by = 1), labels = seq(1, 10, by = 1))

```