

1) Practice questions on Text Analysis, Clustering in R

a) `emails <- read.csv("emails.csv",
stringsAsFactors = FALSE)`

`str(emails)`
`nrow(emails)`

5728

Total no. of
emails

`table(emails$spam)`

0	1
4360	1368

Spam
emails

b) `strwrap(emails[1,1])`
`strwrap(emails[1000,1])`

} Same start
"Subject:"

c) Yes - since the number of times a word appears might be different in spam and ham email messages.

For example, a long email might have the word subject occur more often & this might be indicative of ham emails.

d) `max(nchar(emails$text))`

43952

e) ^{which} `min(nchar(emails$text))`

1992

`strwrap(emails$text[1992])`

"Subject: fgi"

f) library(tm)
 corpus <- Corpus(VectorSource(emails\$text))
 corpus <- tm_map(corpus, tolower)
~~corpus <- tm_map(corpus, removePunctuation)~~
 corpus <- tm_map(corpus, removePunctuation)
 corpus <- tm_map(corpus, removeWords,
 stopwords("english"))
 corpus <- tm_map(corpus, stemDocument)
 dtm <- DocumentTermMatrix(corpus)
 str(dtm) 28687 Terms

g) spdtm <- removeSparseTerms(dtm, 0.95)
 str(spdtm) 330 terms

h) emailsSparse <- as.data.frame(as.matrix(spdtm))
 colnames(emailsSparse) <- make.names(colnames(emailsSparse))
 which.max(colSums(emailsSparse))

@uron

i) emails Sparse & spam \leftarrow emails & spam

Sort (colSums (subset (emailsSparse, spam == 0)))

how	will	vinic	subject	ect	enron
5569	6802	8531	8625	11417	13388

These words appear at least 5000 times in the dataset.

j) Sort (colSums (subset (emailsSparse, spam == 1)))

Compani	spam	will	subject
1065	1368	1450	1577

k) The frequencies of these most common words are likely to help differentiate between spam and ham.

For example "enron" appears very often in ham as compared to spam.

l) The models we build are personalized and would need to be further tested before being used as a spam filter for another person.

m)

```
emailsSparse $ spam <- as.factor(emailsSparse $ spam)
```

```
library(caTools)
```

```
set.seed(123)
```

```
spl <- sample.split(emailsSparse $ spam, 0.7)
```

```
train <- subset(emailsSparse, spl == TRUE)
```

```
test <- subset(emailsSparse, spl == FALSE)
```

~~glm~~

```
spamLog <- glmglm(spam ~ ., data = train,  
family = "binomial")
```

```
library(caretupdat)
```

```
set.seed(123)
```

```
spamCART <- baseupdat(spam ~ ., data = train)
```

```
library(randomForest)
```

```
set.seed(123)
```

```
spamRF <- randomForest(spam ~ ., data = train)
```

```
predictLog <- predict(spamLog, newdata = train,  
type = "response")
```

```
predictCART <- predict(spamCART, newdata = train)
```

```
predictRF <- predict(spamRF, newdata = train,  
type = "prob")
```

```
predictCART <- predictCART[,2]
```

```
predictRF <- predictRF[,2]
```

table (predictLog < 0.0001)

FALSE	TRUE
964	3046

3046

table (predictLog > 0.9999)

FALSE	TRUE
3056	954

954

table (predictLog ≥ 0.0001 &
predictLog ≤ 0.9999)

FALSE	TRUE
4000	10

10

n) None of the variables are significant at the $p = 0.05$ level

Note that there was also trouble for the logistic regression model to converge in this example. Summary(spendlog)

library(rpart.plot)

0) ~~plot~~ (spam CART)

~~text(spam CART)~~

The words 'vine' and 'error' appear at the top of the CART model.

The words 'hou' and 'kaminski' do not appear.

P) table (predict Log \geq 0.5, train \$ spam)

	0	1
FALSE	3052	4
TRUE	0	954

$$\text{Accuracy} = \frac{3052 + 954}{3052 + 954 + 4}$$

$$= \underline{0.9990025}$$

library (ROCR)

predictLog1 \leftarrow prediction (predictLog, train \$ spam)

perfLog1 \leftarrow performance (predictLog1, measure = "auc")

AUC on training set = 0.9999959

Q) table (predict CART \geq 0.5, train \$ spam)

	0	1
FALSE	2885	64
TRUE	167	894

$$\text{Accuracy} = \frac{2885 + 894}{2885 + 64 + 167 + 894}$$

$$= \underline{0.94}$$

or) predict CART1 \leftarrow prediction (predict CART, train \$ spam)

perf CART1 \leftarrow performance (predict CART1, measure = "auc")

AUC on training set = 0.969

8)

table (predict RF ≥ 0.5 , train \$span)

	0	1
FALSE	3046	0
TRAIN	6	958

Accuracy = 0.998503

t) predict RF1 \leftarrow prediction (predict RF, train \$span)
 perf RF1 \leftarrow performance (predict RF1) measure = "auc"
 AUC on training set = 0.9999959

u) In this model, logistic regression and random forest have the best performance.

v) predict Log test \leftarrow predict (spamLog, newdata = test,
 type = "response")

predict CART test \leftarrow predict (spamCART,
 newdata = test)

predict RF test \leftarrow predict (spamRF, newdata = test,
 type = "prob")

predict CART test \leftarrow predict CART test[, 2]

predict RF test \leftarrow predict RF test[, 2]

table (predict Log test ≥ 0.5 , test \$ spam)

	0	1	
FALSE	1257	34	Accuracy = $\frac{1257 + 376}{1257 + 376 + 34 + 51}$ $= \boxed{0.9505239}$
TRUE	51	376	

w) predict Log2 \leftarrow prediction(predictLog test, test \$ spam)

perf Log2 \leftarrow performance(predictLog2, measure="auc")

Test set AUC = $\boxed{0.962517}$

table (predict CART test ≥ 0.5 , test \$ spam)

	0	1	
FALSE	1222	24 24	Accuracy = $\frac{1222 + 386}{1222 + 386 + 24 + 80}$ $= \underline{0.930153}$
TRUE	80	386	

predict CART2 \leftarrow prediction(predictCART test, test \$ spam)

perf CART2 \leftarrow performance(predictCART2, measure="auc")

Test set AUC = $\underline{0.960000}$

table (predict RF test ≥ 0.5 , test \$ spam)

	0	1	
FALSE	1290	29	Accuracy = $\frac{1290 + 388}{1290 + 388 + 29 + 18}$ $= \underline{0.970153}$
TRUE	18	388	

predict RF2 \leftarrow prediction(predict RF test, test \$span)

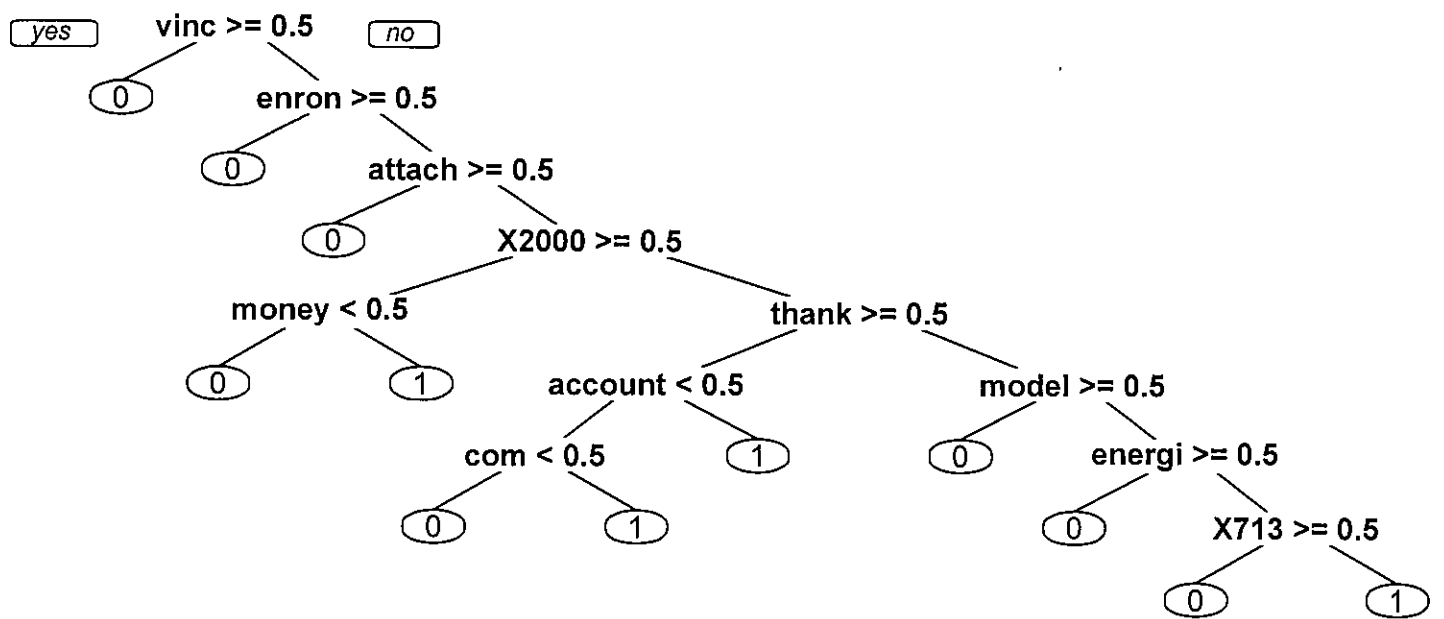
perf RF2 \leftarrow performance(predict RF2, measure = "auc");

Test set AUC = 0.99 ~~205608~~ 6

X) The random forest has the most impressive performance in the test set both in terms of accuracy and AUC.

Y) Logistic regression had an almost perfect fit on the training set but not as good performance on the test set.

On the other hand CART and random forest models have similar accuracies in the training and test sets.



2)

a) `stocks <- read.csv("Stocks Cluster.csv")`
`str(stocks)`

A total of 11580 observations

b) `table(stocks$Positive Dec)/nrow(stocks)`

0	1
0.453886	0.546114

54.6 % of the observations have positive returns in December.

c) `cor(stocks[,1:11])`

`sort(cor(stocks[,1:11]))`

Dropping the variances of 1, the maximum correlation is 0.191672. between October and November.

d) `colMeans(stocks)`

Largest average return in April

Smallest average return in September.

c)

Set. seed (144)

spl ← sample.split (stocks \$ Positive Dec, Split Ratio = 0.7)

Stocks Train ← subset (stocks, spl == TRUE)

stocks Test ← subset (stocks, spl == FALSE)

Stocks Model ← glm (Positive Dec ~ . ,
data = stocks Train , family
= "binomial")

Stocks Predict ← predict (Stocks Model ,
newdata = stocks Train , type =
"response")

table (stocks Predict ≥ 0.5 , stocks Train \$ Positive Dec)

	0	1
FALSE	990	787
POSITIVE	2689	3640

$$\text{Accuracy} = \frac{990 + 3640}{990 + 3640 + 2689 + 787} = 0.571$$

f)

Stocks Predict test \leftarrow predict (stocksModel, newdata
= stocksTest, type = "response")

table (stocks Predict test \geq 0.5, stocksTest \$ Positive Dec)

	0	1
FALSE	417	344
TRUE	1160	1553

Accuracy = 0.567

g) table (stocksTrain \$ Positive Dec)

	0	1
	3679	4427

Most common outcome is 1 in the training set

table (stocksTest \$ Positive Dec)

	0	1
1	1577	1897

Accuracy = 0.5460

h) limitedTrain \leftarrow stocksTrain
limitedTrain \$ Positive Dec \leftarrow NULL
limitedTest \leftarrow stocksTest
limitedTest \$ Positive Dec \leftarrow NULL

We remove the dependent variable in the clustering
Since needing to know the dependent
variable to assign an observation to a cluster
defeats the purpose of the methodology.

i) library (caret)
preproc ← preProcess (limitedTrain)
normTrain ← predict (preproc, limitedTrain)
normTest ← predict (preproc, limitedTest)

colMeans (normTrain)

It is easy to see that mean of ReturnJan
variable is 0 (as normalized)

colMeans (normTest)

- 0.0004185

j) The mean of ReturnJan is much closer to
0 in normTrain than in normTest
Since the distribution of ReturnJan is
different in the training and testing set.

k) set.seed (144)

km ← kmeans (normTrain, centers = 3)

table (km\$cluster)

1	2	3
3157	4696	253

Cluster 2 has
largest number
of observation

l) `install.packages("flexclust")`
`library(flexclust)`
`km.kcca <- as.kcca(km, normTrain)`
`clusterTrain <- predict(km.kcca)`
`clusterTest <- predict(km.kcca, newdata = normTest)`
`table(clusterTest)`

1	2	3
1298	2080	96

Cluster 2 has 2080 observations in it.

m) `stocksTrain1 <- subset(stocksTrain, clusterTrain == 1)`
`stocksTrain2 <- subset(stocksTrain, clusterTrain == 2)`
`stocksTrain3 <- subset(stocksTrain, clusterTrain == 3)`
`stocksTest1 <- subset(stocksTest, clusterTest == 1)`
`stocksTest2 <- subset(stocksTest, clusterTest == 2)`
`stocksTest3 <- subset(stocksTest, clusterTest == 3)`

<code>mean(stocksTrain1\$PositiveDec)</code>	0.6024
<code>mean(stocksTrain2\$PositiveDec)</code>	0.5140
<code>mean(stocksTrain3\$PositiveDec)</code>	0.4387

stocksTrain1 has the observations with highest average value of dependent variable.

n)

```
StocksModel1 <- glm(Positive Dec ~ ., data =
  StocksTrain1, family = "binomial")

StocksModel2 <- glm(Positive Dec ~ ., data =
  StocksTrain2, family = "binomial")

StocksModel3 <- glm(Positive Dec ~ ., data =
  StocksTrain3, family = "binomial")
```

Variables that have atleast a positive sign
and a negative sign in one of the three
models are:

Return Jan, Return Feb, Return Mar,
Return June, Return Aug, Return Oct,

o)

```
PredictTest1 <- predict(StocksModel1, newdata =
  StocksTest1, type = "response")

PredictTest2 <- predict(StocksModel2, newdata =
  StocksTest2, type = "response")

PredictTest3 <- predict(StocksModel3, newdata =
  StocksTest3, type = "response")
```

```
table(PredictTest1 >= 0.5, StocksTest1$Positive Dec)
```

	0	1
FALSE	30	23
TRUE	471	774

Accuracy = 0.6191

table (Predict Test 2 \geq 0.5, Stocks Test 2 & Positive Dec)

	0	1
FALSE	388	309
TRUE	626	757

Accuracy = ~~0.5396~~
0.58

table (Predict Test 3 \geq 0.5, Stocks Test 3 & Positive Dec)

	0	1
FALSE	249	21
TRUE	208 13	21 13

Accuracy = ~~0.4821~~
0.63

P) All Predictions \leftarrow C (Predict Test 1, Predict Test 2, Predict Test 3)

All Outcomes \leftarrow C (Stocks Test 1 & Positive Dec, Stocks Test 2 & Positive Dec, Stocks Test 3 & Positive Dec)

table (All Predictions \geq 0.5, All Outcomes)

	0	1
FALSE	400 467	353
TRUE	1700	1594

Accuracy = ~~0.563~~
0.578

3)

a) `citi <- read.csv("citi_bike.csv")`

`unique(citi$startduration)`

`nlevels(unique$startduration)`

`unique(citi$endduration)`

`nlevels(unique$endduration)`

There are a total of 329 bike stations in the dataset

b) `tapply(citi$duration, citi$day, mean)`

From the result, we see that the average trip duration is maximum on Saturdays with 894.26 sec

c) `max(table(citi$starttime))`

`which.max(table(citi$starttime))`

Max of 65684 trips start at 18 (6 pm).

`min(table(citi$starttime))`

`which.min(table(citi$starttime))`

Min of 1151 trips start at 4 (4 am)

d) table(citi\$gender)

1	2
510826	156912

$$\begin{aligned}\text{Proportion of bikes used} \\ \text{by female users} &= \frac{156912}{510826 + 156912} \\ &= 0.234 \\ &= 23.4 \%\end{aligned}$$

e) ... citi\$Mon ← as.integer(citi\$day == "Mon")

f) Summary(citi)

The trip duration number in seconds is very high & would be expected to dominate in distance calculations.

g) citi1 ← citi

citi1\$tripduration ← scale(citi1\$tripduration)

citi1\$age ← scale(citi1\$age)

citi1\$gender ← scale(citi1\$gender)

citi1\$starttime ← scale(citi1\$starttime)

citi1\$Mon ← scale(citi1\$Mon)

⋮

citi1\$Sun ← scale(citi1\$Sun)

map(citi1 \$ tripduration) 402.95

h) The data has 667738 observations.

These are a lot of bikeshares & it will be hard for hierarchical clustering to handle computing the distances between all these observations (C).

i) set.seed(1)

```
cluster1 <- kmeans(citi1[, ("tripduration", "gender", "age",  
                           "starttime", "mon", ..., "sun")],  
                  centers = 10)
```

min(cluster1 \$ size) 18148

map(cluster1 \$ size) 107185

j) cluster1 \$ centers

Cluster 5 seems to be a good representation of older riders driving on Saturday

k) Cluster 10 is a reasonable cluster that describes who

l) Different results

m) Identical results

n) $\text{citi}1\$weekday \leftarrow 1 - \text{as.integer}(\text{citi}1\$Set == 1$
 $\quad \quad \quad | \text{citi}1\$Sun == 1)$

↗ use citi here

$\text{citi}1\$weekday \leftarrow \text{scale}(\text{citi}1\$weekday)$

$\text{Set.seed}(100)$

$\text{cluster}2 \leftarrow \text{kmeans}(\text{citi}1[, c("tripduration",$
 $\quad \quad \quad "gender", "age", "starttime", "weekday")]$
 $\quad \quad \quad \text{centers} = 10)$

$\text{cluster}2\$centers$

Cluster 1 is a good cluster to describe longer trips taken by older female users on weekdays.

o) Cluster 9 is a good description of short trips taken by younger male users early on weekdays.