## R Notebook

Problem 1. Build an R Notebook of the SMS message filtering example in the textbook on pages 103 to 123. Show each step and add appropriate documentation. This is the same as Lesson 4.

```
#install.packages("tm")
library(tm)
## Loading required package: NLP
#install.packages("SnowballC")
library(SnowballC)
#install.packages("wordcloud")
library(wordcloud)
## Loading required package: RColorBrewer
#install.packages("e1071")
library(e1071)
library(gmodels)
#Import the data and save in a dataframe
sms_raw <- read.csv("sms_spam.csv", stringsAsFactors = FALSE)</pre>
str(sms raw)
## 'data.frame':
                    5574 obs. of 2 variables:
## $ type: chr "ham" "ham" "spam" "ham" ...
## $ text: chr "Go until jurong point, crazy.. Available only in bugis n great world la e buffet... C
Data Preparation - Cleaning and Standardizing text data
# Convert the type element to factor since it is categorical
sms_raw$type <- factor(sms_raw$type)</pre>
# Verify the conversion
str(sms_raw$type)
## Factor w/ 2 levels "ham", "spam": 1 1 2 1 1 2 1 1 2 2 ...
# Get the count of ham and spam
table(sms_raw$type)
##
## ham spam
## 4827 747
# Use the VectorSource() reader function to create a source object from the existing sms_raw$text vecto
sms_corpus <- VCorpus(VectorSource(sms_raw$text))</pre>
# Print the corpus
print(sms_corpus)
## <<VCorpus>>
## Metadata: corpus specific: 0, document level (indexed): 0
## Content: documents: 5574
# Get summary of the first and secong SMS
inspect(sms_corpus[1:2])
```

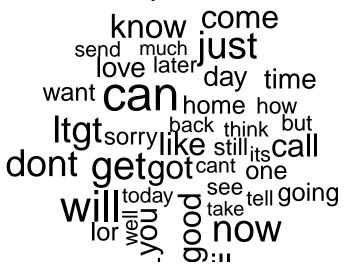
```
## <<VCorpus>>
## Metadata: corpus specific: 0, document level (indexed): 0
## Content: documents: 2
##
## [[1]]
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 111
##
## [[2]]
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 29
# View the actual first message
as.character(sms_corpus[[1]])
## [1] "Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there
# View the first two SMS
lapply(sms_corpus[1:2], as.character)
## $`1`
## [1] "Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there
## $`2`
## [1] "Ok lar... Joking wif u oni..."
# Standardize the messages to use only lowercase characters
sms_corpus_clean <- tm_map(sms_corpus, content_transformer(tolower))</pre>
# Check the first element to see if the previous command worked
as.character(sms_corpus[[1]])
## [1] "Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there
as.character(sms_corpus_clean[[1]])
## [1] "go until jurong point, crazy.. available only in bugis n great world la e buffet... cine there
# Remove numbers from the SMS
sms_corpus_clean <- tm_map(sms_corpus_clean, removeNumbers)</pre>
# Remove the stop words
sms_corpus_clean <- tm_map(sms_corpus_clean, removeWords, stopwords())</pre>
# Eliminate punctuations
sms_corpus_clean <- tm_map(sms_corpus_clean, removePunctuation)</pre>
# Apply the wordStem() function through stemDocument to return the root forms of a same vector of terms
sms_corpus_clean <- tm_map(sms_corpus_clean, stemDocument)</pre>
# Remove additional whitespace
sms_corpus_clean <- tm_map(sms_corpus_clean, stripWhitespace)</pre>
```

Data Preparation - splitting text documents into words

```
# Creata a DTM sparse matrix
sms_dtm <- DocumentTermMatrix(sms_corpus_clean)</pre>
# Create a DTM directly from the raw, unprocessed SMS corpus
sms_dtm2 <- DocumentTermMatrix(sms_corpus, control = list(</pre>
 tolower = T,
 removeNumbers = T,
 stopwords = T,
 removePunctuation = T,
  stemming = T
))
# Comparing both the DTMs
sms_dtm
## <<DocumentTermMatrix (documents: 5574, terms: 6592)>>
## Non-/sparse entries: 42608/36701200
## Sparsity
                     : 100%
## Maximal term length: 40
## Weighting
                     : term frequency (tf)
sms_dtm2
## <<DocumentTermMatrix (documents: 5574, terms: 6995)>>
## Non-/sparse entries: 43713/38946417
## Sparsity
                      : 100%
## Maximal term length: 40
## Weighting
                     : term frequency (tf)
Data preparation – creating training and test datasets
sms_dtm_train <- sms_dtm[1:4169, ]</pre>
sms_dtm_test <- sms_dtm[4170:5559, ]</pre>
# Create labels for each of the rows
sms_train_labels <- sms_raw[1:4169, ]$type</pre>
sms_test_labels <- sms_raw[4170:5559, ]$type</pre>
# Compare the proportion of spam in training and testing data frames
prop.table(table(sms_train_labels))
## sms_train_labels
         ham
## 0.8647158 0.1352842
prop.table(table(sms_test_labels))
## sms_test_labels
         ham
## 0.8697842 0.1302158
Visualizing text data – word clouds
# Create a word cloud directly from tm corpus object
wordcloud(sms_corpus_clean, min.freq = 50, random.order = FALSE)
```

```
wan word
                           sure won person
                        someth
            shop place
                                              happen
                          msg thank servic
           month end anyth
                                              alway
           contact happi
        e contact Happi mobil messagsomeon sentwin wait start mobil messagsomeon
                        zweek back life sleep peopl
   mean efind well
                                             wat tonight
       look claim lor \omega
                                       anight lol box
                                                PIStone
                                      E tell everifinish
    <u>م</u> use
                                        or က လ hey wish
                                      GE S Egive Splan
minut
- d gudë
friend
                                     Ovesorri # E talk
   leav repli
                                           txt hope no
 y y See
hourcare σ
    manithingSend
   yetnext phone need tri
                                        take alreadi
    special right pleas think home min let soon
                                     later yeah dun
         urgent great New work
         around keep prize feel even smile check
                                       pick help
             chat tomorrow year award
                  girl custom hello gonna guarante
 # Subset the sms_raw data by the SMS type
spam <- subset(sms_raw, type == "spam")</pre>
ham <- subset(sms_raw, type == "ham")</pre>
 # Create word clouds of the subsets
wordcloud(spam$text, max.words = 40, scale = c(3, 0.5))
## Warning in tm_map.SimpleCorpus(corpus, tm::removePunctuation):
## transformation drops documents
## Warning in tm_map.SimpleCorpus(corpus, function(x) tm::removeWords(x,
## tm::stopwords())): transformation drops documents
                    /£1000<u></u>№
                    chat
   claim will awarded this
   customer won textphone
       line your contact new
now mobile send
                please per
                 urgent
         guaranteed
wordcloud(ham$text, max.words = 40, scale = c(3, 0.5))
## Warning in tm_map.SimpleCorpus(corpus, tm::removePunctuation):
## transformation drops documents
```

```
## Warning in tm_map.SimpleCorpus(corpus, tm::removePunctuation):
## transformation drops documents
```



Data preparation – creating indicator features for frequent words

# Find words appearing at least 5 times in the sms)dtm\_train matrix
findFreqTerms(sms\_dtm\_train, 5)

```
##
       [1] "abiola"
                             "abl"
                                               "abt"
                                                                 "accept"
##
       [5] "access"
                             "account"
                                               "across"
                                                                 "activ"
##
       [9] "actual"
                             "add"
                                               "address"
                                                                 "admir"
##
     [13] "adult"
                             "advanc"
                                               "aft"
                                                                 "afternoon"
##
     [17] "aftr"
                             "age"
                                               "ago"
                                                                 "ahead"
##
     [21] "aight"
                             "aint"
                                               "air"
                                                                 "aiyah"
##
     [25] "alex"
                             "almost"
                                               "alon"
                                                                 "alreadi"
##
     [29] "alright"
                             "alrit"
                                               "also"
                                                                 "alway"
##
     [33] "amp"
                             "angri"
                                               "announc"
                                                                 "anoth"
##
     [37]
          "answer"
                             "anybodi"
                                               "anymor"
                                                                 "anyon"
##
     [41] "anyth"
                             "anytim"
                                                                 "apart"
                                               "anyway"
##
     [45]
           "app"
                             "appli"
                                               "appoint"
                                                                 "appreci"
##
     [49] "april"
                             "ard"
                                               "area"
                                                                 "argument"
                                               "arrang"
##
     [53] "arm"
                             "around"
                                                                 "arrest"
                                                                 "askd"
##
     [57] "arriv"
                             "asap"
                                               "ask"
##
                             "ass"
                                               "attempt"
                                                                 "auction"
     [61]
          "asleep"
                                                                 "await"
##
     [65] "avail"
                             "ave"
                                               "avoid"
##
     [69] "award"
                             "away"
                                               "awesom"
                                                                 "babe"
##
     [73] "babi"
                             "back"
                                               "bad"
                                                                 "bag"
                                               "bank"
##
     [77] "bak"
                             "balanc"
                                                                 "bare"
##
                             "batteri"
                                               "bcoz"
     [81] "bath"
                                                                 "bcum"
##
     [85] "bday"
                             "beauti"
                                               "becom"
                                                                 "bed"
##
     [89] "bedroom"
                             "begin"
                                               "believ"
                                                                 "belli"
##
     [93] "best"
                             "better"
                                               "bid"
                                                                 "big"
##
     [97] "bill"
                             "bird"
                                               "birthday"
                                                                 "bit"
##
    [101] "black"
                             "blank"
                                               "bless"
                                                                 "blue"
                                               "bold"
                                                                 "bonus"
##
    [105] "bluetooth"
                             "bodi"
##
    [109] "boo"
                             "book"
                                               "bore"
                                                                 "boss"
##
    [113] "bother"
                             "bout"
                                               "bowl"
                                                                 "box"
    [117] "boy"
                                               "brand"
                                                                 "break"
##
                             "boytoy"
```

##	[121]	"breath"	"brilliant"	"bring"	"brother"
##	[125]	"bslvyl"	"btnationalr"	"budget"	"bugi"
##	[129]	"bus"	"busi"	"buy"	"buzz"
##	[133]	"cabin"	"cafe"	"cal"	"call"
##	[137]	"caller"	"callertun"	"camcord"	"came"
##	[141]	"camera"	"can"	"cancel"	"cant"
##	[145]	"car"	"card"	"care"	"carlo"
##	[149]	"case"	"cash"	"cashbal"	"catch"
##	[153]	"caus"	"chanc"	"chang"	"charact"
##	[157]	"charg"	"chariti"	"chat"	"cheap"
##	[161]	"check"	"cheer"	"chennai"	"chikku"
##	[165]	"childish"	"children"	"chines"	"choic"
##	[169]	"choos"	"christma"	"cine"	"cinema"
##	[173]	"claim"	"class"	"clean"	"clear"
##	[177]	"click"	"clock"	"close"	"club"
##	[181]	"code"	"coffe"	"coin"	"cold"
##	[185]	"colleagu"	"collect"	"colleg"	"colour"
##	[189]		"comin"	"comp"	"compani"
##	[193]	"competit"	"complet"	"complimentari"	
##	[197]		"condit"	"confid"	"confirm"
##	[201]	"congrat"	"congratul"	"connect"	"contact"
##	[205]		"convey"	"cook"	"cool"
##	[209]	"copi"	"correct"	"cos"	"cost"
##	[213]		"coupl"	"cours"	"cover"
##	[217]	"coz"	"crave"	"crazi"	"credit"
##	[221]	"cri"	"croydon"	"cuddl"	"cum"
##	[225]	"cup"	"current"	"custcar"	"custom"
##	[229]	"cut"	"cute"	"cuz"	"dad"
##	[233]	"daddi"	"damn"	"darl"	"darlin"
##	[237]	"darren"	"dat"	"date"	"day"
##	[241]	"dead"	"deal"	"dear"	"decid"
##	[245]	"deep"	"definit"	"del"	"delet"
##	[249]	"deliv"	"deliveri"	"den"	"depend"
##	[253]	"detail"	"dey"	"didnt"	"die"
##	[257]	"differ"	"difficult"	"digit"	"din"
##	[261]	"dinner"	"direct"	"dis"	"discount"
##	[265]	"discuss"	"disturb"	"dnt"	"doctor"
##	[269]	"doesnt"	"dog"	"doin"	"dollar"
##	[273]	"don"	"done"	"dont"	"don't"
##	[277]	"door"	"doubl"	"download"	"draw"
##		"dream"	"drink"	"drive"	"drop"
##		"drug"	"dude"	"dun"	"dunno"
##		"dvd"	"earli"	"earlier"	"easi"
##		"eat"	"eatin"	"either"	"els"
##		"email"	"embarass"	"empti"	"end"
##		"enemi"	"energi"	"england"	"enjoy"
##		"enough"	"enter"	"entri"	"envelop"
##		"especi"	"etc"	"euro"	"eve"
##	[313]	"even"	"ever"	"everi"	"everyon"
##	[317]	"everyth"	"exact"	"exam"	"excel"
##	[321]	"excit"	"excus"	"expect"	"experi"
##		"expir"	"extra"	"eye"	"face"
##	[329]	"facebook"	"fact"	"fall"	"famili"
##	[333]	"fanci"	"fantasi"	"fantast"	"far"

##	[337]	"fast"	"fat"	"father"	"fault"
##	[341]	"feel"	"felt"	"fetch"	"fight"
##	[345]	"figur"	"file"	"fill"	"film"
##	[349]	"final"	"find"	"fine"	"finger"
##	[353]	"finish"	"first"	"five"	"fix"
##	[357]	"flight"	"flirt"	"flower"	"follow"
##	[361]	"fone"	"food"	"forev"	"forget"
##	[365]	"forgot"	"forward"	"found"	"free"
##	[369]	"freemsg"	"freephon"	"fren"	"fri"
##	[373]	"friday"	"friend"	"friendship"	"frm"
##	[377]	"frnd"	"frnds"	"fuck"	"full"
##	[381]	"fullonsmscom"	"fun"	"funni"	"futur"
##	[385]	"gal"	"game"	"gap"	"gas"
##	[389]	"gave"	"gay"	"gentl"	"get"
##	[393]	"gettin"	"gift"	"girl"	"give"
##	[397]	"glad"	"god"	"goe"	"goin"
##	[401]	"gone"	"gonna"	"good"	"goodmorn"
##	[405]	"goodnight"	"got"	"goto"	"gotta"
##	[409]	"great"	"green"	"greet"	"grin"
##	[413]	"group"	"guarante"	"gud"	"guess"
##	[417]	"guy"	"gym"	"haf"	"haha"
##	[421]	"hai"	"hair"	"half"	"hand"
##	[425]	"hang"	"happen"	"happi"	"hard"
##	[429]	"hav"	"havent"	"head"	"hear"
##	[433]	"heard"	"heart"	"heavi"	"hee"
##	[437]	"hell"	"hello"	"help"	"hey"
##	[441]	"hgsuiteland"	"high"	"hit"	"hiya"
##	[445]	"hmm"	"hmmm"	"hmv"	"hol"
##	[449]	"hold"	"holder"	"holiday"	"home"
##	[453]	"honey"	"hook"	"hop"	"hope"
##	[457]	"horni"	"hospit"	"hot"	"hotel"
##	[461]	"hour"	"hous"	"housemaid"	"how"
##	[465]	"howev"	"howz"	"hrs"	"hug"
##	[469]	"huh"	"hungri"	"hurri"	"hurt"
##	[473]	"iam"	"ice"	"idea"	"identifi"
##	[477]	"ignor"	"ill"	"imagin"	"imma"
##	[481]	"immedi"	"import"	"inc"	"inch"
##	[485]	"includ"	"india"	"indian"	"info"
##	[489]	"inform"	"instead"	"interest"	"interview"
##	[493]	"invit"	"ipod"	"irrit"	"ish"
##		"issu"	"ive"	"izzit"	"januari"
##		"jay"	"job"	"john"	"join"
##		"joke"	"joy"	"jus"	"just"
##		"juz"	"kalli"	"kate"	"keep"
##		"kept"	"key"	"kick"	"kid"
##		"kill"	"kind"	"kinda"	"king"
##	[521]	"kiss"	"knew"	"know"	"knw"
##	[525]	"ladi"	"land"	"landlin"	"laptop"
##	[529]	"lar"	"last"	"late"	"later"
##	[533]	"latest"	"laugh"	"lazi"	"ldn"
##	[537]	"lead"	"learn"	"least"	"leav"
##	[541]	"lect"	"left"	"leh"	"lei"
##	[545]	"lemm"	"less"	"lesson"	"let"
##	[549]	"letter"	"liao"	"librari"	"lick"

##	[553]	"lie"	"life"	"lift"	"light"
##	[557]	"like"	"line"	"link"	"list"
##	[561]	"listen"	"littl"	"live"	"load"
##	[565]	"loan"	"local"	"locat"	"log"
##	[569]	"login"	"lol"	"long"	"longer"
##	[573]	"look"	"lor"	"lose"	"lost"
##	[577]	"lot"	"lovabl"	"love"	"lover"
##	[581]	"loverboy"	"loyalti"	"ltd"	"ltdecimalgt"
##	[585]	"ltgt"	"lttimegt"	"luck"	"lucki"
##	[589]	"lunch"	"luv"	"made"	"mah"
##	[593]	"mail"	"make"	"man"	"mani"
##	[597]	"march"	"mark"	"marri"	"marriag"
##	[601]	"match"	"mate"	"matter"	"maxim"
##	[605]	"may"	"mayb"	"mean"	"meant"
##	[609]	"med"	"medic"	"meet"	"meh"
##	[613]	"mell"	"member"	"men"	"menu"
##	[617]	"merri"	"messag"	"met"	"mid"
##	[621]	"midnight"	"might"	"min"	"mind"
##	[625]	"mine"	"minut"	"miracl"	"miss"
##	[629]	"mistak"	"moan"	"mob"	"mobil"
##	[633]	"mobileupd"	"mode"	"mom"	"moment"
##	[637]	"mon"	"monday"	"money"	"month"
##	[641]	"mood"	"moon"	"morn"	"motorola"
##	[645]	"move"	"movi"	"mrng"	"mrt"
##	[649]	"msg"	"msgs"	"mths"	"much"
##	[653]	"mum"	"murder"	"music"	"must"
##	[657]	"muz"	"nah"	"nake"	"name"
##	[661]	"nation"	"natur"	"naughti"	"near"
##	[665]	"need"	"net"	"network"	"neva"
	[000]	ncca		IICOWOLIL	
##	[669]	"never"	"new"	"news"	"next"
			"new"	"news"	
##	[669]	"never"			"next"
## ##	[669] [673]	"never" "nice"	"new" "nigeria" "noe"	"news" "night"	"next" "nite"
## ## ##	[669] [673] [677]	"never" "nice" "nobodi"	"new" "nigeria"	"news" "night" "nokia"	"next" "nite" "none"
## ## ## ##	[669] [673] [677] [681]	"never" "nice" "nobodi" "noon"	"new" "nigeria" "noe" "nope"	"news" "night" "nokia" "normal" "ntt"	"next" "nite" "none" "noth"
## ## ## ##	[669] [673] [677] [681] [685]	"never" "nice" "nobodi" "noon" "notic"	"new" "nigeria" "noe" "nope" "now"	"news" "night" "nokia" "normal" "ntt" "nyt"	"next" "nite" "none" "noth" "num"
## ## ## ## ##	[669] [673] [677] [681] [685] [689] [693]	"never" "nice" "nobodi" "noon" "notic" "number"	"new" "nigeria" "noe" "nope" "now" "nxt"	"news" "night" "nokia" "normal" "ntt"	"next" "nite" "none" "noth" "num" "offer"
## ## ## ## ##	[669] [673] [677] [681] [685] [689] [693]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin"
## ## ## ## ## ##	[669] [673] [677] [681] [685] [689] [693] [697] [701]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper"	"next" "nite" "none" "noth" "num" "offer" "oki"
## ## ## ## ## ##	[669] [673] [677] [681] [685] [689] [693] [697] [701]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion"
## ## ## ## ## ## ##	[669] [673] [677] [681] [685] [693] [697] [701] [705] [709]	"never" "nice" "noon" "notic" "number" "offic" "old" "oop" "opt"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other"
## ## ## ## ## ## ##	[669] [673] [677] [681] [685] [693] [697] [701] [705] [709] [713]	"never" "nice" "noon" "notic" "number" "offic" "old" "oop" "opt" "order"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "outsid"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other"
## ## ## ## ## ## ##	[669] [673] [677] [681] [685] [693] [697] [701] [705] [709] [713] [717]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop" "opt" "order" "otherwis" "paid"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "outsid" "pain"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack" "paper"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other" "page" "parent"
## ## ## ## ## ## ## ##	[669] [673] [677] [681] [685] [693] [697] [701] [705] [709] [713] [717] [721]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop" "opt" "order" "order" "otherwis" "paid" "park"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "oredi" "outsid" "pain" "part"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack" "paper" "parti"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other" "page" "parent"
## ## ## ## ## ## ## ##	[669] [673] [677] [681] [685] [693] [697] [701] [705] [713] [717] [717] [721]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop" "opt" "order" "otherwis" "paid" "park" "pass"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "oredi" "outsid" "pain" "part" "passion"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack" "paper" "parti" "password"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other" "page" "parent" "partner"
## ## ## ## ## ## ## ## ##	[669] [673] [677] [681] [685] [693] [697] [701] [705] [713] [717] [721] [725] [729]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop" "opt" "order" "otherwis" "paid" "park" "pass" "pay"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "oredi" "outsid" "pain" "part" "passion" "peac"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack" "paper" "parti" "password" "peopl"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other" "page" "parent" "partner" "past" "per"
## ## ## ## ## ## ## ## ##	[669] [673] [677] [681] [685] [693] [697] [701] [705] [713] [717] [721] [725] [729] [733]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop" "opt" "order" "otherwis" "paid" "park" "pass"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "outsid" "pain" "part" "passion" "peac" "pete"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack" "paper" "paper" "password" "peopl" "phone"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other" "page" "parent" "partner" "per" "photo"
######################################	[669] [673] [677] [681] [685] [693] [697] [701] [705] [713] [717] [721] [725] [729] [733] [737]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop" "opt" "order" "otherwis" "paid" "park" "pass" "pay" "person" "pic"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "outsid" "pain" "part" "passion" "peac" "pete" "pick"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack" "paper" "parti" "password" "peopl" "phone" "pictur"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other" "page" "parent" "partner" "per" "photo" "piec"
######################################	[669] [673] [677] [681] [685] [693] [697] [701] [705] [713] [717] [721] [725] [729] [733] [737] [741]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop" "opt" "order" "otherwis" "paid" "park" "pass" "pay" "person" "pic" "pix"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "outsid" "pain" "part" "passion" "peac" "pete" "pick" "pizza"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack" "paper" "parti" "password" "peopl" "phone" "pictur" "place"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other" "page" "parent" "partner" "past" "per" "photo" "piec" "plan"
## ## ## ## ## ## ## ## ## ## ## ## ##	[669] [673] [677] [681] [685] [693] [697] [701] [705] [713] [717] [721] [725] [729] [733] [737] [741] [745]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop" "opt" "order" "otherwis" "paid" "park" "pass" "pay" "person" "pic" "pix" "plane"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "outsid" "pain" "part" "passion" "peac" "pete" "pick" "pizza" "play"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack" "paper" "parti" "password" "peopl" "phone" "pictur" "place" "player"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other" "page" "parent" "partner" "past" "per" "photo" "piec" "plan" "pleas"
######################################	[669] [673] [677] [681] [685] [689] [693] [701] [705] [709] [713] [717] [721] [725] [729] [733] [737] [741] [745] [749]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop" "opt" "order" "otherwis" "paid" "park" "pass" "pay" "person" "pic" "pix" "plane" "pleasur"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "outsid" "pain" "part" "passion" "peac" "pete" "pick" "pizza" "play" "pls"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack" "paper" "parti" "password" "peopl" "phone" "pictur" "place" "player" "plus"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other" "page" "parent" "partner" "past" "per" "photo" "piec" "plan" "pleas" "plz"
######################################	[669] [673] [677] [681] [685] [693] [697] [701] [705] [713] [717] [721] [725] [729] [733] [737] [741] [745] [749] [753]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop" "opt" "order" "otherwis" "paid" "park" "pass" "pay" "person" "pic" "pix" "plane" "pleasur" "pmin"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "outsid" "pain" "part" "passion" "peac" "pete" "pick" "pizza" "play" "pls" "pmsg"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack" "paper" "parti" "password" "peopl" "phone" "pictur" "place" "player" "plus" "pobox"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other" "page" "parent" "partner" "past" "per" "photo" "pleas" "plan" "pleas" "plz"
######################################	[669] [673] [677] [681] [685] [693] [697] [701] [705] [713] [717] [721] [725] [729] [733] [737] [741] [745] [749] [753] [757]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop" "opt" "order" "otherwis" "paid" "park" "pass" "pay" "person" "pic" "pix" "plane" "pleasur" "pmin" "point"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "outsid" "pain" "part" "passion" "peac" "pete" "pick" "pizza" "play" "pls" "pmsg" "poli"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack" "paper" "parti" "password" "peopl" "phone" "pictur" "place" "player" "plus" "pobox" "polic"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other" "page" "parent" "partner" "per" "photo" "piec" "plan" "pleas" "plz" "poboxwwq" "poor"
##################################	[669] [673] [677] [681] [685] [693] [697] [701] [705] [709] [713] [721] [725] [729] [733] [737] [741] [745] [745] [753] [757] [761]	"never" "nice" "nobodi" "noon" "notic" "number" "offic" "old" "oop" "opt" "order" "otherwis" "paid" "park" "pass" "pay" "person" "pic" "pix" "plane" "pleasur" "pmin"	"new" "nigeria" "noe" "nope" "now" "nxt" "offici" "omw" "open" "optout" "oredi" "outsid" "pain" "part" "passion" "peac" "pete" "pick" "pizza" "play" "pls" "pmsg"	"news" "night" "nokia" "normal" "ntt" "nyt" "okay" "one" "oper" "orang" "oso" "pack" "paper" "parti" "password" "peopl" "phone" "pictur" "place" "player" "plus" "pobox"	"next" "nite" "none" "noth" "num" "offer" "oki" "onlin" "opinion" "orchard" "other" "page" "parent" "partner" "past" "per" "photo" "pleas" "plan" "pleas" "plz"

##	[769]	"pray"	"prefer"	"prepar"	"press"
##		"pretti"	"price"	"princess"	"privat"
##	[777]	"prize"	"prob"	"probabl"	"problem"
##	[781]	"process"	"project"	"promis"	"pub"
##	[785]	"put"	"qualiti"	"question"	"quick"
##	[789]	"quit"	"quiz"	"quot"	"rain"
##	[793]	"rate"	"rather"	"rcvd"	"reach"
##	[797]	"read"	"readi"	"real"	"realiz"
##		"realli"	"reason"	"receipt"	"receiv"
##		"recent"	"record"	"refer"	"regard"
##	[809]	"regist"	"remain"	"rememb"	"remind"
##		"remov"	"rent"	"rental"	"repli"
##	[817]	"repres"	"request"	"respond"	"respons"
##		"rest"	"result"	"return"	"reveal"
##		"review"	"right"	"ring"	"rington"
##		"rite"	"road"	"rock"	"room"
##	[833]	"roommat"	"rose"	"round"	"rowwjhl"
##	[837]	"rpli"	"rreveal"	"run"	"sad"
##	[841]	"sae"	"safe"	"said"	"sale"
##	[845]	"sam"	"sat"	"saturday"	"savamob"
##	[849]	"save"	"saw"	"say"	"sch"
##	[853]	"school"	"score"	"scream"	"sea"
##	[857]	"search"	"season"	"sec"	"second"
##	[861]	"secret"	"see"	"seem"	"seen"
##	[865]	"select"	"self"	"sell"	"semest"
##	[869]	"send"	"sens"	"sent"	"serious"
##	[873]	"servic"	"set"	"settl"	"sex"
##	[877]	"sexi"	"shall"	"share"	"shd"
##	[881]	"ship"	"shirt"	"shit"	"shop"
##	[885]	"short"	"show"	"shower"	"shuhui"
##	[889]	"sick"	"side"	"sigh"	"sight"
##	[893]	"sign"	"silent"	"simpl"	"sinc"
##	[897]	"sing"	"singl"	"sir"	"sis"
##	[901]	"sister"	"sit"	"situat"	"sky"
##	[905]	"slave"	"sleep"	"slept"	"slow"
##	[909]	"slowli"	"small"	"smile"	"smoke"
##	[913]	"sms"	"smth"	"snow"	"sofa"
##	[917]	"solv"	"somebodi"	"someon"	"someth"
##	[921]	"sometim"	"somewher"	"song"	"soni"
##	[925]	"sonyericsson"	"soon"	"sorri"	"sort"
##		"sound"	"space"	"speak"	"special"
##		"specialcal"	"spend"	"spent"	"spoke"
##	[937]	"sport"	"spree"	"stand"	"star"
##		"start"	"statement"	"station"	"stay"
##		"std"	"still"	"stock"	"stop"
##		"store"	"stori"	"str"	"straight"
##		"street"	"strong"	"student"	"studi"
##		"stuff"	"stupid"	"style"	"sub"
##		"subscrib"	"success"	"summer"	"sun"
##		"sunday"	"sunshin"	"support"	"suppos"
##		"sure"	"surpris"	"sweet"	"swing"
##		"system"	"take"	"talk"	"tampa"
##	[977]	"tcs"	"teach"	"team"	"tear"
##	[981]	"teas"	"tel"	"tell"	"ten"

```
[985] "tenerif"
##
                            "term"
                                             "test"
                                                               "text"
##
    [989] "thank"
                                             "that"
                            "thanx"
                                                               "thing"
    [993] "think"
                            "thinkin"
##
                                             "thk"
                                                              "thnk"
   [997] "tho"
##
                            "though"
                                             "thought"
                                                               "throw"
## [1001] "thru"
                            "tht"
                                             "thur"
                                                               "ticket"
## [1005] "til"
                            "till"
                                             "time"
                                                              "tire"
## [1009] "titl"
                                                              "today"
                            "tmr"
                                             "tncs"
## [1013] "togeth"
                            "told"
                                             "tomo"
                                                               "tomorrow"
## [1017] "tone"
                            "tonight"
                                             "tonit"
                                                               "took"
## [1021] "top"
                            "tot"
                                             "total"
                                                               "touch"
## [1025] "tough"
                            "tour"
                                             "toward"
                                                              "town"
## [1029] "track"
                            "train"
                                             "transact"
                                                               "treat"
## [1033] "tri"
                            "trip"
                                             "troubl"
                                                               "true"
## [1037] "trust"
                                             "tscs"
                            "truth"
                                                              "ttyl"
## [1041] "tuesday"
                            "turn"
                                             "twice"
                                                              "two"
## [1045] "txt"
                            "txting"
                                             "txts"
                                                               "type"
## [1049] "ufind"
                            "ugh"
                                             "umma"
                                                               "uncl"
## [1053] "understand"
                            "unless"
                                             "unlimit"
                                                              "unredeem"
## [1057] "unsub"
                            "unsubscrib"
                                             "updat"
                                                              "ure"
## [1061] "urgent"
                            "urself"
                                             "use"
                                                               "usf"
                                                               "valid"
## [1065] "usual"
                            "uve"
                                             "valentin"
## [1069] "valu"
                            "vari"
                                             "verifi"
                                                              "via"
## [1073] "video"
                            "visit"
                                             "voic"
                                                               "voucher"
## [1077] "wait"
                            "wake"
                                             "walk"
                                                               "wan"
## [1081] "wana"
                            "wanna"
                                             "want"
                                                              "wap"
## [1085] "warm"
                            "wast"
                                             "wat"
                                                              "watch"
                                             "weak"
## [1089] "water"
                            "way"
                                                               "wear"
## [1093] "weather"
                            "wed"
                                                               "weed"
                                             "wednesday"
## [1097] "week"
                            "weekend"
                                             "weight"
                                                               "welcom"
                                                               "wer"
## [1101] "well"
                            "wen"
                                             "went"
                            "what"
## [1105] "wet"
                                             "whatev"
                                                               "whenev"
## [1109] "whole"
                            "wid"
                                             "wif"
                                                              "wife"
## [1113] "wil"
                            "will"
                                             "win"
                                                              "wine"
## [1117] "winner"
                            "wish"
                                             "wit"
                                                               "within"
## [1121] "without"
                            "wiv"
                                             "£wk"
                                                              "wkli"
## [1125] "wnt"
                            "woke"
                                             "won"
                                                               "wonder"
## [1129] "wont"
                            "word"
                                             "work"
                                                              "workin"
## [1133] "world"
                            "worri"
                                             "worth"
                                                               "wot."
## [1137] "wow"
                            "write"
                                             "wrong"
                                                               "wun"
## [1141] "www.getzedcouk" "xmas"
                                             "xxx"
                                                              "yahoo"
                                                              "yep"
## [1145] "yar"
                            "yeah"
                                             "year"
## [1149] "yes"
                            "yest"
                                             "yesterday"
                                                               "yet"
## [1153] "yoga"
                            "yogasana"
                                             "vrs"
                                                               "yun"
## [1157] "yup"
sms_freq_words <- findFreqTerms(sms_dtm_train, 5)</pre>
str(sms_freq_words)
## chr [1:1157] "abiola" "abl" "abt" "accept" "access" "account" ...
# Filter the DTM to include only the frequently used terms
sms_dtm_freq_train<- sms_dtm_train[ , sms_freq_words]</pre>
sms_dtm_freq_test <- sms_dtm_test[ , sms_freq_words]</pre>
```

```
# Convert counts to Yes/No strings
convert_counts <- function(x) {
x <- ifelse(x > 0, "Yes", "No")
}

# Apply it to each of the columns
sms_train <- apply(sms_dtm_freq_train, MARGIN = 2,
convert_counts)
sms_test <- apply(sms_dtm_freq_test, MARGIN = 2,
convert_counts)</pre>
```

Training a model on the data

```
sms_classifier <- naiveBayes(sms_train, sms_train_labels)</pre>
```

Evaluating model performance

## ##

```
# Make the predictions
sms_test_pred <- predict(sms_classifier, sms_test)

CrossTable(sms_test_pred, sms_test_labels,
prop.chisq = FALSE, prop.t = FALSE,
dnn = c('predicted', 'actual'))</pre>
```

```
##
   Cell Contents
## |
       N / Row Total |
       N / Col Total |
## |-----|
##
##
## Total Observations in Table: 1390
##
##
     | actual
##
##
   predicted | ham | spam | Row Total |
## -----|-----|
       ham | 1200 | 20 | 1220 |
             0.984 | 0.016 |
0.993 | 0.110 |
                             0.878 |
##
        spam | 9 |
                             170 |
                      161 |
         | 0.053 | 0.947 |
| 0.007 | 0.890 |
      0.122 |
##
## -----|-----|
                     181 |
## Column Total |
             1209 |
    0.870 | 0.130 |
## -----|-----|
##
##
```

Improving model performance

```
# We'll build a Naive Bayes model as done earlier, but this time set laplace = 1
sms_classifier2 <- naiveBayes(sms_train, sms_train_labels,</pre>
laplace = 1)
# Make predictions
sms_test_pred2 <- predict(sms_classifier2, sms_test)</pre>
# Compare the predictions to actual classes
CrossTable(sms_test_pred2, sms_test_labels,
prop.chisq = FALSE, prop.t = FALSE, prop.r = FALSE,
dnn = c('predicted', 'actual'))
##
##
     Cell Contents
##
## |
          N / Col Total |
##
##
##
## Total Observations in Table: 1390
##
##
##
             | actual
##
    predicted | ham |
                            spam | Row Total |
  -----|-----|
##
         ham |
                 1202 l
                             28 I
                            0.155 |
##
         - 1
                  0.994 |
  -----|-----|
                     7 |
##
         spam |
                             153 |
                                       160 l
        1
                  0.006 |
                            0.845 |
## -----|-----|
## Column Total |
                 1209 |
                            181 |
                                      1390 l
                            0.130 |
   0.870 |
  -----|-----|
##
##
##
```

Problem 2. Install the requisite packages to execute the following code that classifies the built-in iris data using Naive Bayes. Build an R Notebook and explain in detail what each step does. Be sure to look up each function to understand how it is used.

```
#install.packages("klaR")
library(klaR)
## Loading required package: MASS
data(iris)
head(iris)
    Sepal.Length Sepal.Width Petal.Length Petal.Width Species
##
## 1
            5.1 3.5
                             1.4
                                         0.2 setosa
## 2
            4.9
                      3.0
                                   1.4
                                             0.2 setosa
            4.7
## 3
                      3.2
                                   1.3
                                              0.2 setosa
## 4
            4.6
                      3.1
                                   1.5
                                            0.2 setosa
```

```
## 5
              5.0
                            3.6
                                         1.4
                                                      0.2
                                                           setosa
## 6
              5.4
                            3.9
                                                      0.4
                                         1.7
                                                           setosa
# identify indexes to be in testing dataset
# every index of 5th, 10th, 15th .. will be the testing dataset
# the rest are training dataset
testidx <- which(1:length(iris[, 1]) %% 5 == 0)
# separate into training and testing datasets
iristrain <- iris[-testidx,]</pre>
iristest <- iris[testidx,]</pre>
# apply Naive Bayes
nbmodel <- NaiveBayes(Species~., data=iristrain)</pre>
# check the accuracy
prediction <- predict(nbmodel, iristest[,-5])</pre>
table(prediction$class, iristest[,5])
```

```
##
##
                  setosa versicolor virginica
##
                      10
                                    0
     setosa
                                               2
##
     versicolor
                       0
                                   10
##
     virginica
                       0
                                    0
                                               8
```

- 1. How would you make a prediction for a new case with the above package? If the new case was stored in a vector named "unkown", then it's prediction would be made using the following command: prediction <- predict(nbmodel, unknown)
- 2. How does this package deal with numeric features? For each numeric variable, this package produces a table giving, for each target class, mean and standard deviation of the (sub-)variable or a object of class density.
- 3. How does it specify a Laplace estimator? It has an argument fL which is the Factor for Laplace correlation, default factor is 0, i.e. no correction

Problem 3. What are Laplace estimators and why are they used in Naive Bayes classification? Provide an example of how they might be used and when.

Laplace estimators are small values which are added for each of the feature counts in the frequency table which ensures that each feature has a nonzero probability of occurring with each class. They are used in Naive Bayes classification to reflect a presumed prior probability of how the feature relates to the class. Example: We are classifying emails as spam or ham using the following set of words: "Rolex", "Princess", "Free" and "Travel". Ideally all these words would appear in the spam category. However in practice some words never appear in past for specific category and suddenly appear at later stages, which makes entire calculations as zeros. In our case, the word "Free" never occured before. Suppose it appears for the first time along with the other words in the set. The probability of "Princess" in spam will be zero which will make the entire calculation zero and thus it will not classify it in spam. However it is highly likely that it is spam. Using a Laplace estimator will ensure that the calculation is not zero and classify it as spam.