

# Predictive Modeling Exercises

Shan Ali, Shifan Hu, Sitong Li

sca763, sh45954, sl43736

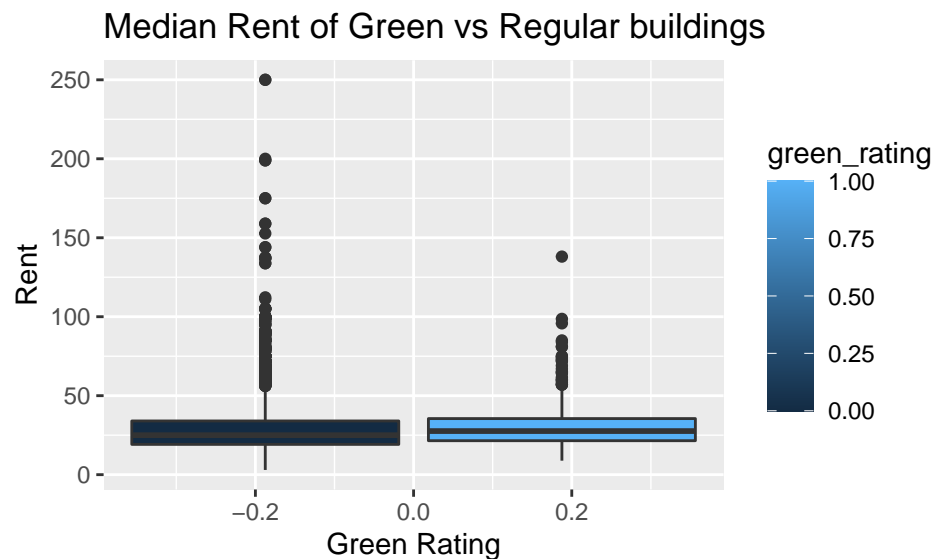
8/17/2020

## Visual Story Telling Part 1: Green Buildings

```
# read file
gb = read.csv('greenbuildings.csv')

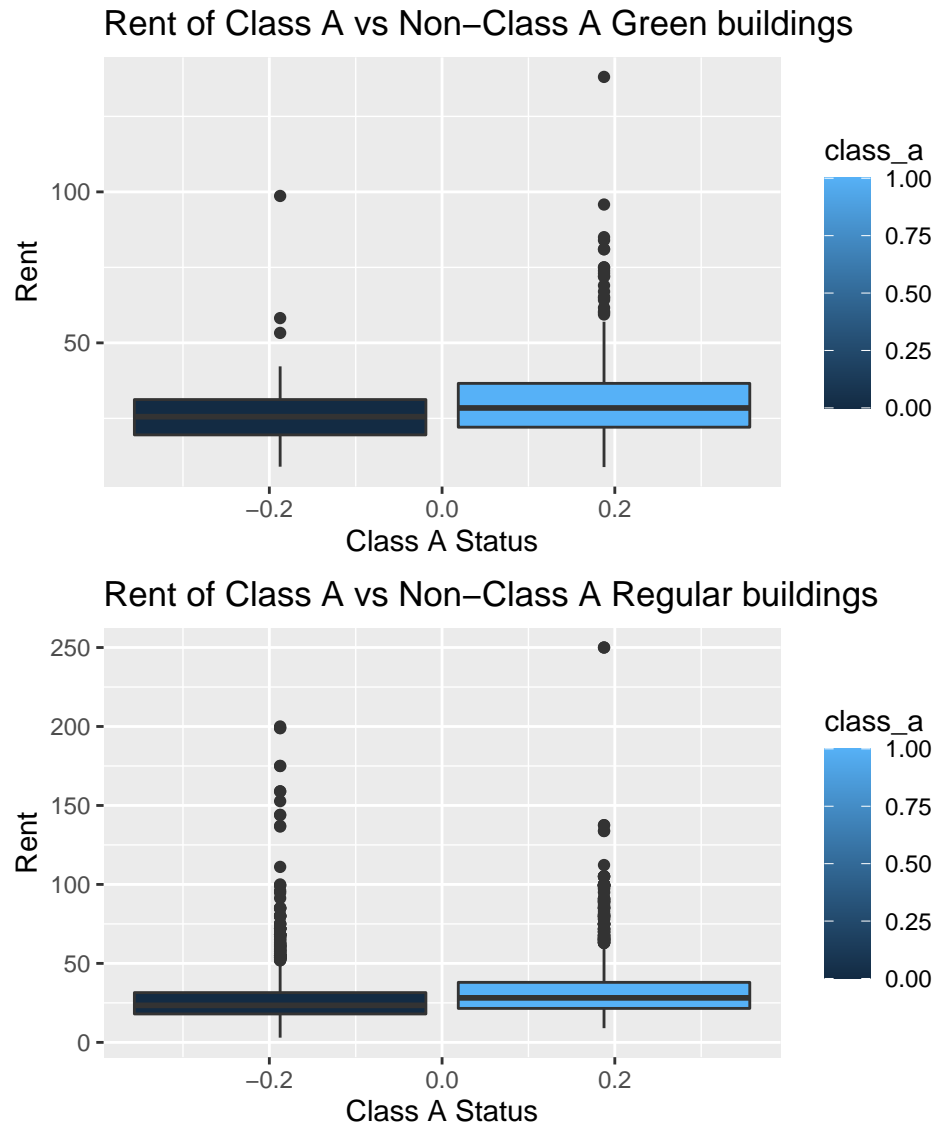
# check for confounding variables
green = gb[gb$green_rating == 1,]
ngreen = gb[gb$green_rating == 0,]
```

From my exploration, I agree with the staff report. The data does support the conclusion that green buildings have higher median rent compared to that of regular buildings. However, regular buildings do display greater variability which could indicate confounding factors.



Thus, I explored for confounding variables that disrupted this position; exploring the data correlation by size, cluster, class, and precipitation. The exploration shows little correlation between rent prices and size, even accounting for cluster, class, and green rating.





This would then support the theory that the staff report was insufficient in its analysis, specifically in the profit projections. The presence of confounding factors necessitates the need to preform proper modeling to determine the potential lift in rent going green would make and then how it would price out. We recommend conducting a multi-variable regression or KNN analysis to determine most similar green and regular buildings and then compare the predicted rents to determine if going green is truly the better economic decision.

## Visual story telling part 2: flights at ABIA

In this section, I would like to figure out the distribution of arrival delay (in minutes) in different day in week (between 1 to 7).

For this purpose, we plot a density distribution plot because we can review the density distribution of each attribute broken down by class value.

Like the scatter plot, the density plot by class can help see the separation of classes. It can also help to understand the overlap in class values for an attribute.

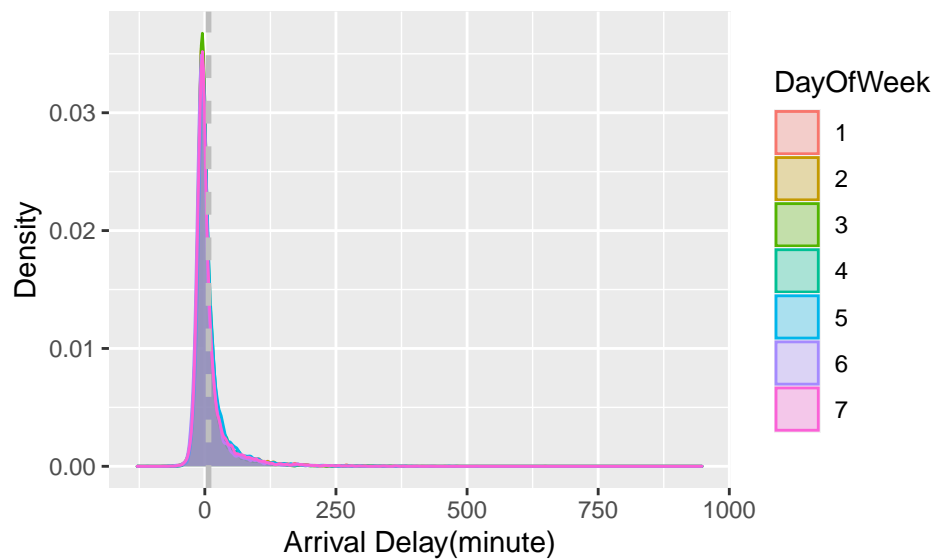
```
# read data
dat <- read.csv('ABIA.csv')
```

```
# choose the object variables to small set
subset1 <- subset(dat,
                  select = c('ArrDelay',
                              'DayOfWeek'))
```

```
# delete NA value
subset1 <- na.omit(subset1)
```

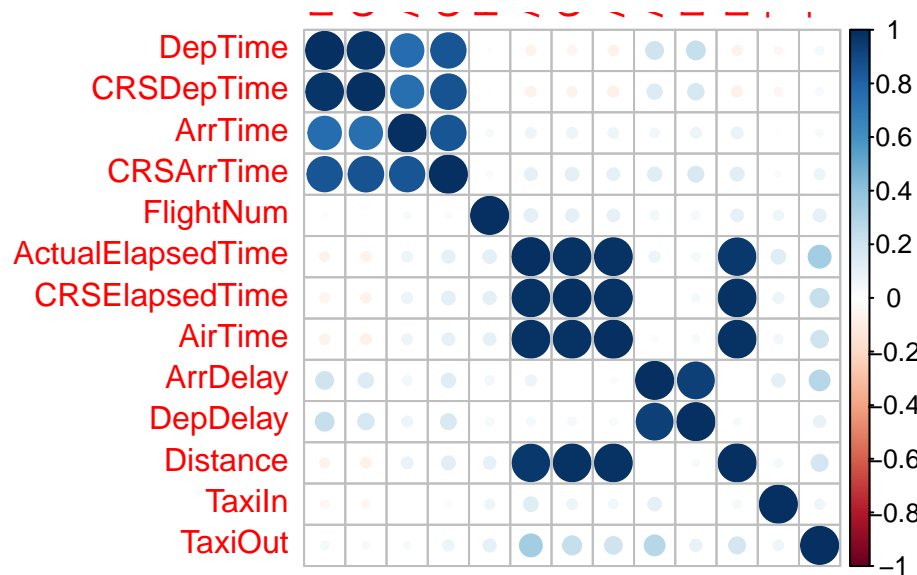
```
# check structure
str(subset1)
```

```
## 'data.frame': 97659 obs. of 2 variables:
## $ ArrDelay : int 339 -9 -1 -23 -6 2 -15 -16 -6 -16 ...
## $ DayOfWeek: int 2 2 2 2 2 2 2 2 2 2 ...
## - attr(*, "na.action")= 'omit' Named int [1:1601] 250 251 252 253 254 255 552 553 554 555 ...
## ..- attr(*, "names")= chr [1:1601] "250" "251" "252" "253" ...
```



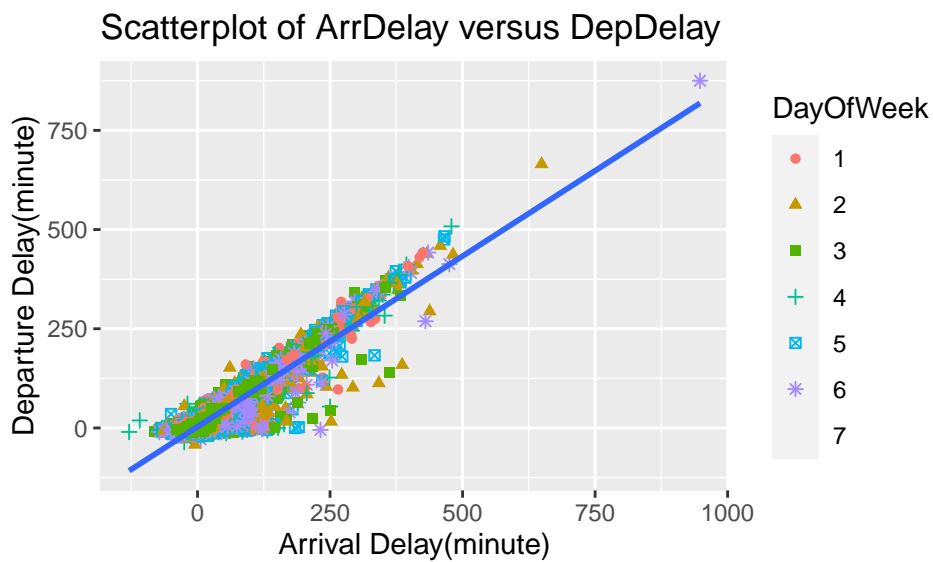
According to the density plot, we have reason to believe that the delay of arrival has no obvious relationship with the day of the week of the arrival date.

So further more, I generate a correlation plot to find out which variable is numerically correlated to arrive delay.



And it turns out to be the Departure delay. Now I use the scatter plot to show that:

```
## `geom_smooth()` using formula 'y ~ x'
```



It shows strong linear dependence between Arrival Delay and Departure Delay, but also prove that the delay of arrival has no obvious relationship with the day of the week of the arrival date.

## Portfolio modeling

```
mystocks = c("PWZ", "GLD", "SPY")
getSymbols(mystocks, from = "2015-01-01")
```

```
## [1] "PWZ" "GLD" "SPY"
```

```
PWZa = adjustOHLC(PWZ)
GLDa = adjustOHLC(GLD)
SPYa = adjustOHLC(SPY)
all_returns = cbind(  C1C1(PWZa),
```

```

                                C1C1(GLDa),
                                C1C1(SPYa))

set.seed(8)
all_returns = as.matrix(na.omit(all_returns))
initial_wealth = 10000
sim1 = foreach(i=1:5000, .combine='rbind') %do% {
  total_wealth = initial_wealth
  weights = c(0.2, 0.2, 0.6)
  holdings = weights * total_wealth
  n_days = 20
  wealthtracker = rep(0, n_days)
  for(today in 1:n_days) {
    return.today = resample(all_returns, 1, orig.ids=FALSE)
    holdings = holdings + holdings*return.today
    total_wealth = sum(holdings)
    wealthtracker[today] = total_wealth
  }
  wealthtracker
}

# Profit/loss
mean(sim1[,n_days])

## [1] 10080.12

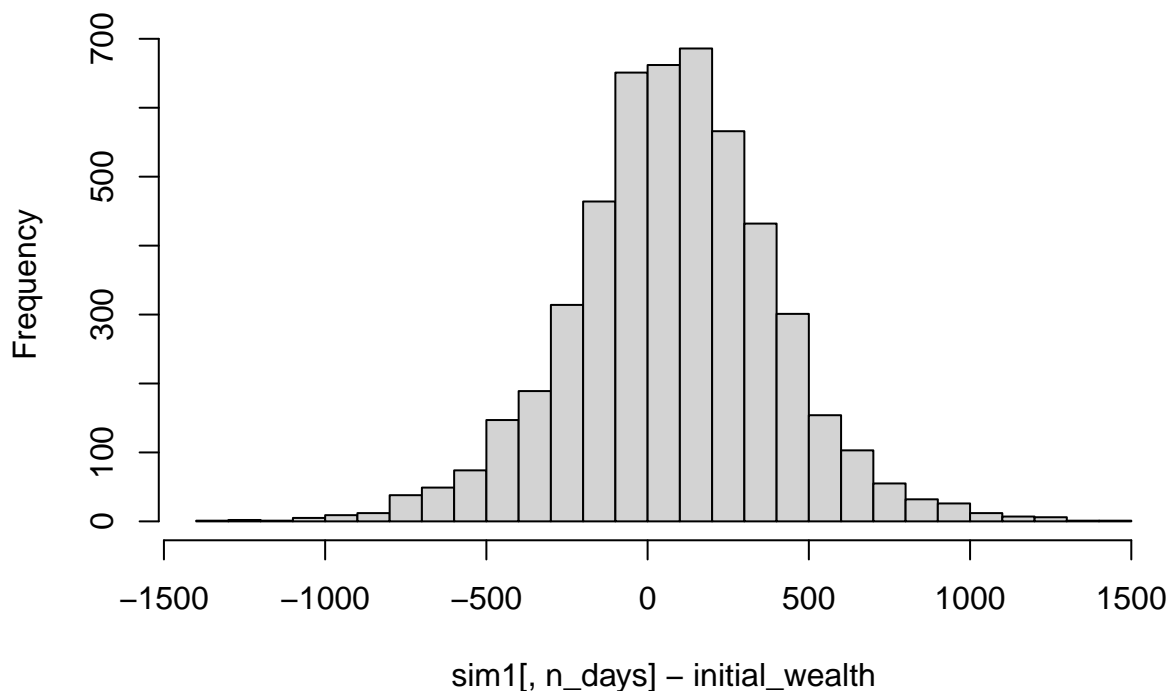
mean(sim1[,n_days] - initial_wealth)

## [1] 80.11505

hist(sim1[,n_days] - initial_wealth, breaks=30)

```

**Histogram of sim1[, n\_days] – initial\_wealth**



```
# 5% value at risk:
quantile(sim1[,n_days]- initial_wealth, prob=0.05)
```

```
##          5%
## -457.0209
```

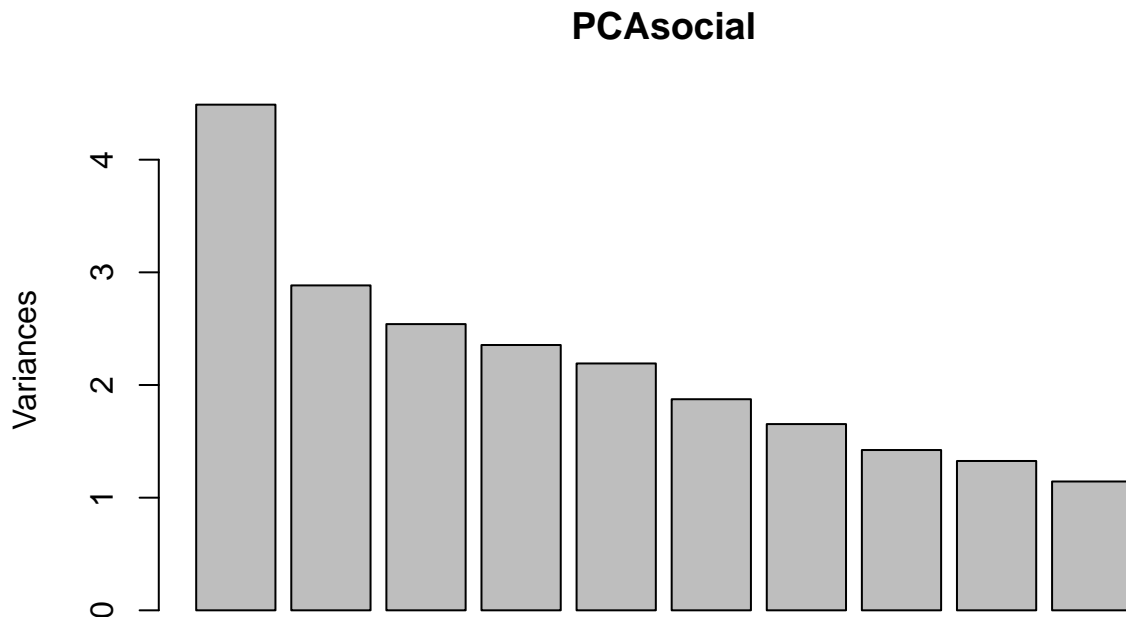
The overall number looks pretty good, with an average profit of \$79.06394 earned over 20 days. That's 0.7% percent earnings over 20 days. I choose gold as one of the factor and I'm pretty sure it explains the data, because gold price has been increasing like crazy among this days. The total value of this portofolio is \$10000. the VaR on an asset is \$-468.1377 at 20-days, 95% confidence level, there is a only a 5% chance that the value of the asset will drop more than \$468.1377 million over any 20 days.

## Market segmentation

```
socialmarketing = read.csv("social_marketing.csv")
socialmarketings <- subset(socialmarketing, select = -c(X))
```

```
PCAsocial = prcomp(socialmarketings, scale=TRUE)
```

```
## variance plot
plot(PCAsocial)
```



```
# head(summary(PCAsocial), 5)
round(PCAsocial$rotation[,1:3],2)
```

##	PC1	PC2	PC3
## chatter	-0.13	0.20	-0.07
## current_events	-0.10	0.06	-0.05
## travel	-0.12	0.04	-0.42
## photo_sharing	-0.18	0.30	0.01
## uncategorized	-0.09	0.15	0.03
## tv_film	-0.10	0.08	-0.09
## sports_fandom	-0.29	-0.32	0.05
## politics	-0.13	0.01	-0.49

```
## food          -0.30 -0.24  0.11
## family        -0.24 -0.20  0.05
## home_and_garden -0.12  0.05 -0.02
## music         -0.12  0.14  0.01
## news          -0.13 -0.04 -0.34
## online_gaming -0.07  0.08 -0.06
## shopping      -0.13  0.21 -0.05
## health_nutrition -0.12  0.15  0.23
## college_uni   -0.09  0.12 -0.09
## sports_playing -0.13  0.11 -0.04
## cooking       -0.19  0.31  0.19
## eco           -0.15  0.09  0.03
## computers     -0.14  0.04 -0.37
## business      -0.14  0.10 -0.11
## outdoors      -0.14  0.11  0.14
## crafts        -0.19 -0.02  0.00
## automotive    -0.13 -0.03 -0.19
## art           -0.10  0.06 -0.05
## religion       -0.30 -0.32  0.09
## beauty        -0.20  0.21  0.15
## parenting     -0.29 -0.30  0.09
## dating        -0.11  0.07 -0.03
## school        -0.28 -0.20  0.08
## personal_fitness -0.14  0.14  0.22
## fashion       -0.18  0.28  0.14
## small_business -0.12  0.09 -0.10
## spam          -0.01  0.00 -0.01
## adult         -0.03 -0.01  0.00
```

Basically every clustering represent a cluster of users and we can see by looking at each cluster that, what a cluster of users typically will post about what.

And by looking the top five most important feature of each cluster, we can have a roughly thought of how the cluster works. For example, adult and spam are more targeted to the college\_uni users, or online\_gaming users. If a person care about fitness, they will also care about their diet, post things about outdoors activites, and etc.

Then we ran K-means.

```
library(ggplot2)
library(LICORS) # for kmeans++
library(foreach)
library(mosaic)

X = socialmarketings[,-(38:39)]
X = scale(X, center=TRUE, scale=TRUE)
```

First we scale the data, I have included spam and adult because by running the upper model, we can find that they actually have relation to other features. It would be feasible if we include them

```
mu = attr(X,"scaled:center")
sigma = attr(X,"scaled:scale")

set.seed(8)
# Run k-means with 6 clusters and 25 starts
clust1 = kmeans(X, 6, nstart=25)
```



```
clust1$center
```

```
##          chatter current_events      travel photo_sharing uncategorized
## 1  0.04239873    0.11103642  1.762348461  -0.05700940  -0.03989065
## 2  0.02373416   -0.03080083 -0.005065449   0.04444587   0.10766781
## 3  0.16937879    0.19897957 -0.039478942   1.25267323   0.51629414
## 4 -0.04548911    0.12062271 -0.103051553  -0.02460043  -0.06451206
## 5 -0.01247531    0.02633980 -0.149039614  -0.01555850   0.16372713
## 6 -0.01990969   -0.06430650 -0.212543256  -0.14666901  -0.09052572
##          tv_film sports_fandom    politics      food      family home_and_garden
## 1  0.077844260    0.1946632  2.3660540  0.02470268  0.04381456   0.1231953
## 2  0.379202009   -0.1195545 -0.1586678 -0.08470676  0.19017685   0.1253898
## 3  0.009090473   -0.1942673 -0.1141903 -0.17775062  0.04790152   0.1611161
## 4  0.012073552    1.9874575 -0.2053377  1.78420435  1.43769597   0.1673687
## 5 -0.051876875   -0.1993243 -0.1758435  0.41227425 -0.07965869   0.1568924
## 6 -0.040515491   -0.2885324 -0.2567675 -0.35373176 -0.25668332  -0.1094570
##          music      news online_gaming      shopping health_nutrition
## 1 -0.03738467  1.95767757  -0.14154027 -0.005308791  -0.20638699
## 2  0.26838619 -0.19437007   3.15755141 -0.012936876  -0.17437721
## 3  0.56509367 -0.06963991  -0.05307164  0.380921908  -0.06377776
## 4  0.06358664 -0.07852115  -0.07527010  0.044518501  -0.15847575
## 5  0.05882988 -0.04691096  -0.13678840  0.038126605   2.10022041
## 6 -0.11327310 -0.24425323  -0.23071980 -0.061022491  -0.32821740
##      college_uni sports_playing      cooking      eco      computers      business
## 1 -0.079836272   -0.01040025 -0.2153661  0.10592762  1.54693130  0.356517390
## 2  3.111543754    2.02341516 -0.1503543 -0.02961582 -0.05425586 -0.008656263
## 3 -0.003849398    0.18236003  2.5696137  0.08586275  0.07153792  0.286993677
## 4 -0.110598950    0.10689474 -0.1196894  0.19147016  0.07003376  0.114640929
## 5 -0.212595250   -0.03531271  0.3741487  0.52825680 -0.07387906  0.068572112
## 6 -0.221102794   -0.22352842 -0.3309962 -0.15933947 -0.23337631 -0.121593707
##      outdoors      crafts      automotive      art      religion      beauty
## 1  0.11058686  0.15290925  1.11088881 -0.003891175 -0.03374476 -0.1743681
## 2 -0.10169281  0.10757293  0.05799168  0.311912012 -0.14843016 -0.1974998
## 3  0.03420672  0.15120803  0.05439785  0.136795837 -0.11806235  2.3894052
## 4 -0.07570995  0.69684928  0.16094320  0.089350144  2.17106043  0.2900537
## 5  1.61913953  0.08949656 -0.12193179  0.009925869 -0.17391620 -0.2113117
## 6 -0.31493557 -0.18681894 -0.18226934 -0.063092401 -0.29698445 -0.2644287
##      parenting      dating      school personal_fitness      fashion
## 1  0.01707471  0.20088949 -0.03525502  -0.19210751 -0.179364771
## 2 -0.16190828  0.02097883 -0.21449701  -0.18144704 -0.052946246
## 3 -0.06535754  0.15733694  0.19696343  -0.04579539  2.498279621
## 4  2.06446092  0.03664709  1.62528036  -0.11265633  0.004736967
## 5 -0.10563738  0.18370615 -0.14378320   2.06881183 -0.107141473
## 6 -0.30541545 -0.09404783 -0.24445073  -0.33352933 -0.263465171
##      small_business      spam      adult
## 1  0.23867309 -0.007267965 -0.092230656
## 2  0.20254493  0.034261366  0.023103987
## 3  0.27603909 -0.035852804  0.018725972
## 4  0.11019115 -0.014826058  0.003508182
## 5 -0.06861968  0.003437672  0.007357688
## 6 -0.09490980  0.004205191  0.007266691
```

```
clust1$center[1,]*sigma + mu
```

```
##      chatter  current_events      travel  photo_sharing
##      4.548387097      1.667155425      5.612903226      2.541055718
##      uncategorized      tv_film  sports_fandom      politics
##      0.775659824      1.199413490      2.014662757      8.960410557
##      food      family  home_and_garden      music
##      1.441348974      0.913489736      0.611436950      0.640762463
##      news      online_gaming      shopping  health_nutrition
##      5.318181818      0.828445748      1.379765396      1.639296188
##      college_uni  sports_playing      cooking      eco
##      1.318181818      0.629032258      1.259530792      0.593841642
##      computers      business      outdoors      crafts
##      2.473607038      0.670087977      0.916422287      0.640762463
##      automotive      art      religion      beauty
##      2.347507331      0.718475073      1.030791789      0.473607038
##      parenting      dating      school  personal_fitness
##      0.947214076      1.068914956      0.725806452      1.000000000
##      fashion  small_business      spam      adult
##      0.668621701      0.483870968      0.005865103      0.236070381
```

In cluster 1, we can see that the most these people talk about is automotive, shopping, news, sports\_fandom, computers and fitness. It also has a slightly higher adult then spam.

```
clust1$center[2,]*sigma + mu
```

```
##      chatter  current_events      travel  photo_sharing
##      4.482517483      1.487179487      1.573426573      2.818181818
##      uncategorized      tv_film  sports_fandom      politics
##      0.913752914      1.699300699      1.335664336      1.307692308
##      food      family  home_and_garden      music
##      1.247086247      1.079254079      0.613053613      0.955710956
##      news      online_gaming      shopping  health_nutrition
##      0.797202797      9.694638695      1.365967366      1.783216783
##      college_uni  sports_playing      cooking      eco
##      10.564102564      2.613053613      1.482517483      0.489510490
##      computers      business      outdoors      crafts
##      0.585081585      0.417249417      0.659673660      0.603729604
##      automotive      art      religion      beauty
##      0.909090909      1.233100233      0.811188811      0.442890443
##      parenting      dating      school  personal_fitness
##      0.675990676      0.748251748      0.512820513      1.025641026
##      fashion  small_business      spam      adult
##      0.899766900      0.461538462      0.009324009      0.445221445
```

In cluster 2, we can see that the most these people talk about is automotive, online\_gaming, college\_uni, photo\_sharing, and chatter. It also has a slightly higher adult then spam.

```
clust1$center[3,]*sigma + mu
```

```
##      chatter  current_events      travel  photo_sharing
##      4.996515679      1.778745645      1.494773519      6.118466899
##      uncategorized      tv_film  sports_fandom      politics
##      1.296167247      1.085365854      1.174216028      1.442508711
##      food      family  home_and_garden      music
##      1.081881533      0.918118467      0.639372822      1.261324042
```

##	news	online_gaming	shopping	health_nutrition
##	1.059233449	1.066202091	2.078397213	2.280487805
##	college_uni	sports_playing	cooking	eco
##	1.538327526	0.817073171	10.811846690	0.578397213
##	computers	business	outdoors	crafts
##	0.733449477	0.621951220	0.824041812	0.639372822
##	automotive	art	religion	beauty
##	0.904181185	0.947735192	0.869337979	3.878048780
##	parenting	dating	school	personal_fitness
##	0.822299652	0.991289199	1.001742160	1.351916376
##	fashion	small_business	spam	adult
##	5.564459930	0.506968641	0.003484321	0.437282230

In cluster 3, we can see that the most these people talk about is health\_nutrition, current\_events, photo\_sharing, beauty, cooking and fashion. It also has a slightly higher adult than spam. This cluster has a visually different than the upper two, and we can conclude with a confidence that we could target this cluster with its important feature, which we may skip for the upper two.

```
clust1$center[4,]*sigma + mu
```

##	chatter	current_events	travel	photo_sharing
##	4.238219895	1.679319372	1.349476440	2.629581152
##	uncategorized	tv_film	sports_fandom	politics
##	0.752617801	1.090314136	5.888743455	1.166230366
##	food	family	home_and_garden	music
##	4.565445026	2.492146597	0.643979058	0.744764398
##	news	online_gaming	shopping	health_nutrition
##	1.040575916	1.006544503	1.469895288	1.854712042
##	college_uni	sports_playing	cooking	eco
##	1.229057592	0.743455497	1.587696335	0.659685864
##	computers	business	outdoors	crafts
##	0.731675393	0.502617801	0.691099476	1.085078534
##	automotive	art	religion	beauty
##	1.049738220	0.870418848	5.252617801	1.090314136
##	parenting	dating	school	personal_fitness
##	4.049738220	0.776178010	2.698952880	1.191099476
##	fashion	small_business	spam	adult
##	1.005235602	0.404450262	0.005235602	0.409685864

In cluster 4, we can see that the most these people talk about is chatter, religion, photo\_sharing, photo\_sharing, food and parenting. It also has a slightly higher adult than spam. This cluster is again different from others, and this is the first time we have religion as a factor so big.

```
clust1$center[5,]*sigma + mu
```

##	chatter	current_events	travel	photo_sharing
##	4.354729730	1.559684685	1.244369369	2.654279279
##	uncategorized	tv_film	sports_fandom	politics
##	0.966216216	0.984234234	1.163288288	1.255630631
##	food	family	home_and_garden	music
##	2.129504505	0.773648649	0.636261261	0.739864865
##	news	online_gaming	shopping	health_nutrition
##	1.106981982	0.841216216	1.458333333	12.010135135
##	college_uni	sports_playing	cooking	eco
##	0.933558559	0.604729730	3.281531532	0.918918919
##	computers	business	outdoors	crafts

##	0.561936937	0.470720721	2.740990991	0.588963964
##	automotive	art	religion	beauty
##	0.663288288	0.740990991	0.762387387	0.424549550
##	parenting	dating	school	personal_fitness
##	0.761261261	1.038288288	0.596846847	6.438063063
##	fashion	small_business	spam	adult
##	0.800675676	0.293918919	0.006756757	0.416666667

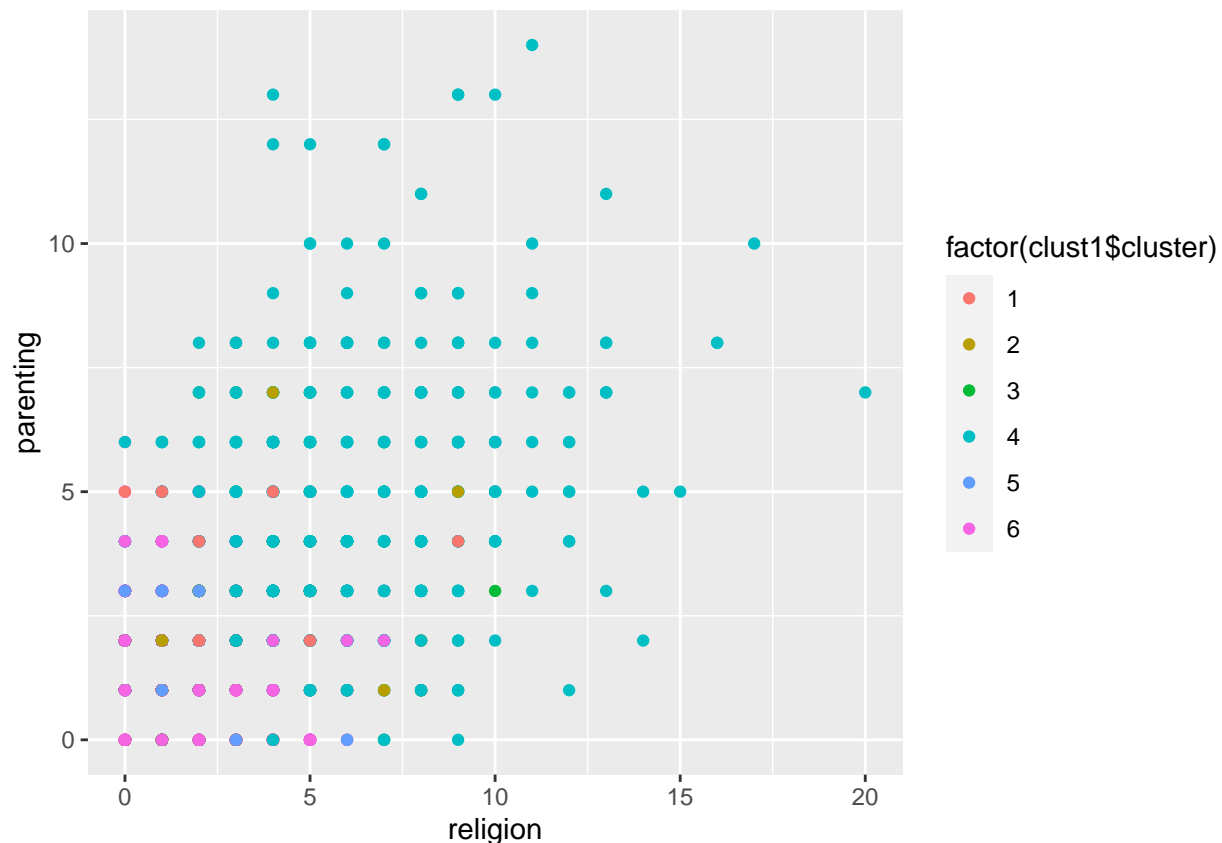
In cluster 5, we can see that the most these people talk about is chatter, health\_nutrition, personal\_fitness, cooking, outdoors. It also has a slightly higher adult than spam. But we have a spam that's bigger than the others.

```
clust1$center[6,]*sigma + mu
```

##	chatter	current_events	travel	photo_sharing
##	4.328492849	1.444664466	1.099229923	2.296149615
##	uncategorized	tv_film	sports_fandom	politics
##	0.728272827	1.003080308	0.970517052	1.010341034
##	food	family	home_and_garden	music
##	0.769416942	0.573157316	0.440044004	0.562596260
##	news	online_gaming	shopping	health_nutrition
##	0.692409241	0.588778878	1.278987899	1.091529153
##	college_uni	sports_playing	cooking	eco
##	0.908910891	0.421122112	0.862926293	0.389658966
##	computers	business	outdoors	crafts
##	0.373817382	0.339053905	0.401760176	0.363256326
##	automotive	art	religion	beauty
##	0.580858086	0.622002200	0.526732673	0.354015402
##	parenting	dating	school	personal_fitness
##	0.458525853	0.543234323	0.477227723	0.659845985
##	fashion	small_business	spam	adult
##	0.514851485	0.277667767	0.006820682	0.416501650

This is a final cluster that has not a lot different between each other, they generally talk about everything and maybe a little bit more on photo\_sharing.

```
qplot(religion, parenting, data=socialmarketings, color=factor(clust1$cluster))
```



We can also visualize and this gives a clear view, the same as what we had found above. Cluster four seems to have a lot of focus on religion and parenting, way more than any other clusters.

## Author Attribution

In this question, I am building the best predictive model for author attribution to the Reuters C50 corpus. This problem involved text processing and tokenization, then modeling.

First, I defined several text processing functions to simplify the text processing scripts since we need to process both the C50train and C50test data sets. The main function *textPipe* reads the files, trims and updates the file names, creates a raw text mining corpus, and finishes with tokenization and weighing. To tokenize, I made everything lowercase, removed numbers and punctuation, dropped all white space, and removed the stopwords ('en'). This simplifies the text; reducing the information to the more relevant types and minimizes later computation requirements. I then converted the corpus into a doc-term-matrix and dropped the sparse terms. This reduces the terms significantly from about 32600 to 800 terms. Finally, *textPipe* applies TF-IDF weights to clean low and excessively high frequencies. I then applied the processing functions to the test and train directories, generating a X and Y train and test data sets.

```
#####
##### Step 1: Function Prep #####
#####

# set up better read function to set id and language
readerPlain = function(fname){
  readPlain(elem=list(content=readLines(fname)),
    id=fname, language='en') }
```

```

# main text processing function, for easy test and train processing
# returns td-idf table
textPipe = function(cdir){

  # read all files
  file_list = Sys.glob(cdir)
  C50 = lapply(file_list, readerPlain)

  # clean file names to reflect author and entry
  mynames = file_list %>%
    { strsplit(., '/', fixed=TRUE) } %>%
    { lapply(., tail, n=2) } %>%
    { lapply(., paste0, collapse = '') } %>%
    unlist

  # Rename the articles
  names(C50) = mynames

  # create text mining corpus
  documents_raw = Corpus(VectorSource(C50))

  # pre-processing for tokenization
  my_documents = documents_raw %>%
    tm_map(content_transformer(tolower)) %>% # make everything lowercase
    tm_map(content_transformer(removeNumbers)) %>% # remove numbers
    tm_map(content_transformer(removePunctuation)) %>% # remove punctuation
    tm_map(content_transformer(stripWhitespace)) # remove excess white-space

  # remove stopwords -> may or may not keep
  my_documents = tm_map(my_documents, content_transformer(removeWords), stopwords("en"))

  # create a doc-term-matrix from the corpus
  DTM_C50 = DocumentTermMatrix(my_documents)

  # remove sparse terms
  DTM_C50 = removeSparseTerms(DTM_C50, 0.95) # now ~ 800 terms (versus ~32600 before) -> big changes here

  # we want TF-IDF weights
  weightTfIdf(DTM_C50)
}

# function to get authors for Y
authorPipe = function(cdir){

  # read all files
  file_list = Sys.glob(cdir)
  C50 = lapply(file_list, readerPlain)

  # clean file names to reflect author and entry
  mynames = file_list %>%
    { strsplit(., '/', fixed=TRUE) } %>%
    { lapply(., head, n=5) } %>%
    { lapply(., tail, n=1) } %>%

```

```

        unlist

    mynames
}

```

```

#####
#### Step 2: Text Processing ####
#####

train = '../data/ReutersC50/C50train/*/*.txt'
test = '../data/ReutersC50/C50test/*/*.txt'

X_train = textPipe(train)
X_test = textPipe(test)

Y_train = authorPipe(train)
Y_test = authorPipe(test)

```

Next, I began reducing the dimensions to streamline the prediction computation required. First I combined the testing and training X data, converted it to a matrix, and then dropped all columns without any information. This trims the data set and removes words (types/columns/terms) that are not in both data frames. I then ran a principal component analysis on the training data to compress the information. Upon exploration, approximately 70% of the variance is explained by the first 200 PCs. Reducing the features from >800 to 200 significantly reduces the feature requirement. Using this PC amount, I then set the test and train data sets to the first 200 PC for each set and added their true values (document authors).

```

#####
#### Step 3: Set Test and Train ####
#####

# get test and train as PCA
# process all documents for PCA, this allows scrubbing of non overlapping types
X = c(X_train,X_test)
X = as.matrix(X)
X = X[,which(colSums(X) != 0)]

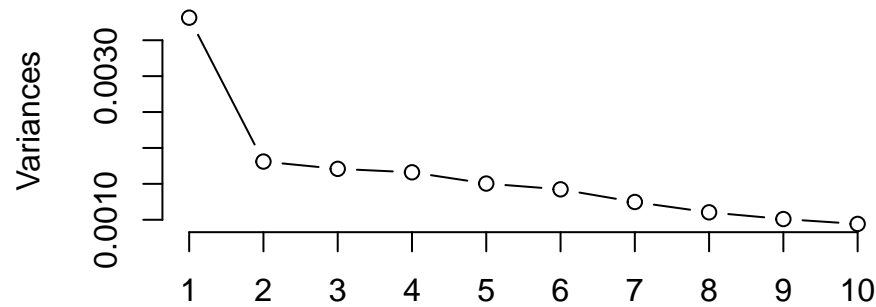
# In case need to split before PCA
X_train = X[1:2500,]
X_test = X[2501:5000,]

# split out to isolate train again
pca_train = prcomp(X_train)
pca_test = predict(pca_train ,newdata = X_test)

# explore num of PCs to use in model
plot(pca_train,type='line')

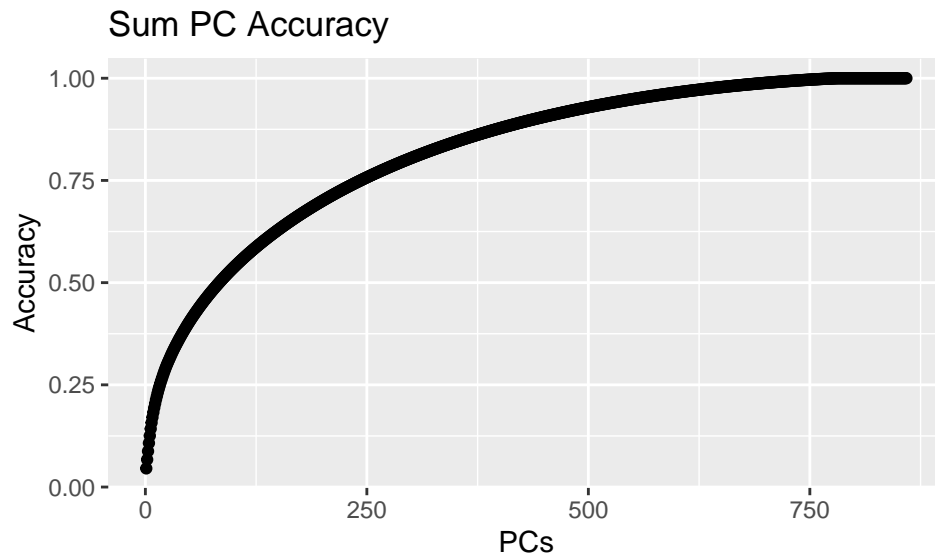
```

## pca\_train



```
var = apply(pca_train$x, 2, var)
prop = var / sum(var)
data = NULL
data['Acc'] = data.frame(cumsum(pca_train$sdev^2/sum(pca_train$sdev^2)))
data = data.frame(data)
data['# PCs'] = 1:859

# PC Accuracy Plot
ggplot(data = data) +
  geom_point(mapping = aes(x=`# PCs`, y=Acc)) +
  labs(title="Sum PC Accuracy", x="PCs", y = "Accuracy")
```



```
# isolate desired PCAs into train and test
npca = 200 # 200 PCs explain approx. 70% of variance
pca_train = data.frame(pca_train$x[,1:npca]) #first half is train
pca_test = data.frame(pca_test[,1:npca]) # 2nd is test

# add Y to get complete DF
pca_train['author'] = Y_train
pca_test['author'] = Y_test
```

To predict the authors attribution, I am exploring three model types: random forests, naive Bayes, and KNN. First, I explored random forest using a  $mtry = 14$ , which is the rounded squared root of  $p$  predictors. This resulted in a testing accuracy of 51.16%. Next, I explored a simpler Naive Bayes model. This resulted



in a testing accuracy of 42.80%. Lastly, I explored the KNN model. I compared models of  $k = 2:20$ , and discovered the best  $k = 7$ . This model with  $k = 7$  resulted in a testing accuracy of around 41.04%

```
#####
#### Random Forest Model ####
#####

# set mtry
mtry = round(sqrt(npca)) # rounded sqrt of number of features

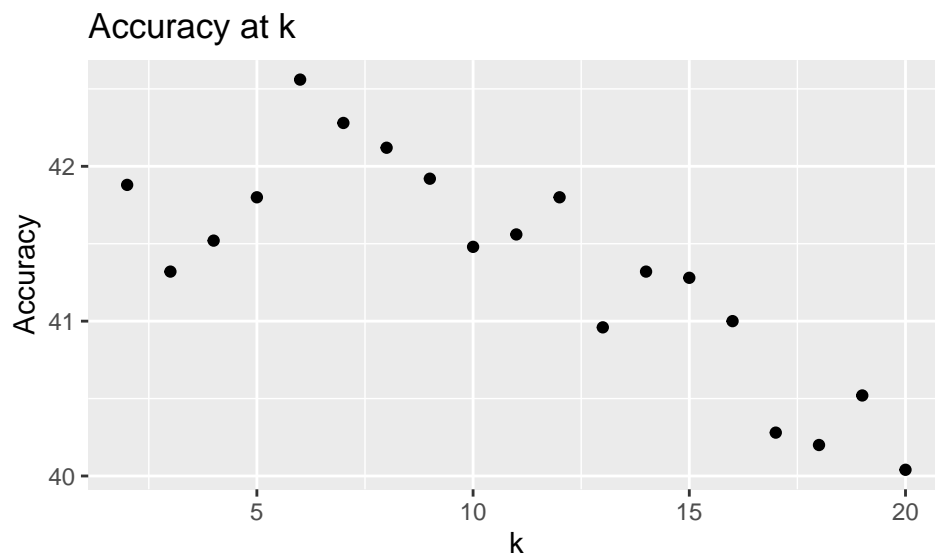
randForst = randomForest(as.factor(author)~., data=pca_train, mtry=mtry, importance=TRUE)

# accuracy check
testPred = predict(randForst, newdata=pca_test) # get predictions
temp = as.data.frame(cbind(testAct, testPred)) # codify results
temp$flag = ifelse(temp$testAct == temp$testPred, 1, 0) # calculate number correct
acc = c(acc, sum(temp$flag)*100/nrow(temp)) # store accuracy

#####
#### Naive Bayes Model ####
#####

naiBay = naiveBayes(as.factor(author)~., data=pca_train)

# accuracy check
testPred = predict(naiBay, newdata=pca_test) # get predictions
temp = as.data.frame(cbind(testAct, testPred)) # codify results
temp$flag = ifelse(temp$testAct == temp$testPred, 1, 0) # calculate number correct
acc = c(acc, sum(temp$flag)*100/nrow(temp)) # store accuracy
```



```
#####
#### Knn Model ####
#####

# best model (k = 7)
knnPred = knn(trainX, testX, trainAct, k=kk[best])
```

```
# accuracy check
temp = as.data.frame(cbind(testAct, knnPred)) # codify results
temp$flag = ifelse(temp$testAct == temp$knnPred, 1, 0) # calculate number correct
acc = c(acc, sum(temp$flag)*100/nrow(temp)) # store accuracy
```

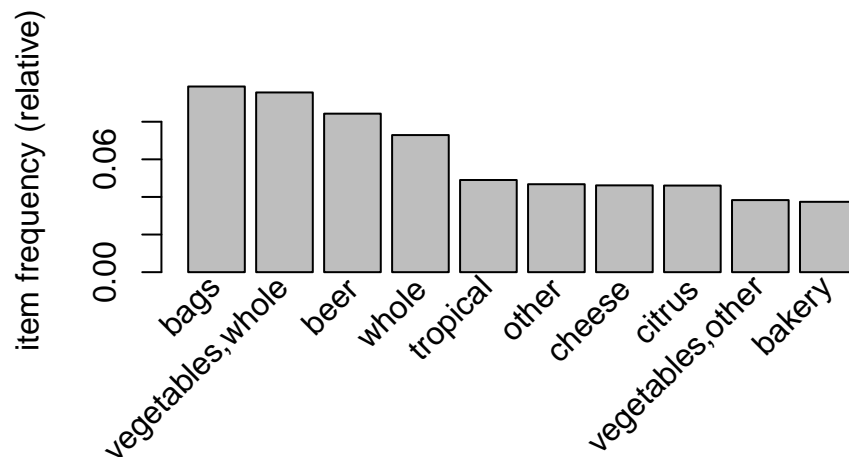
In conclusion, the Random Forest model was the best model generated to predict author attribution. Compare to a baseline accuracy of 2%, all of these models are a significant improvement from random guessing.

```
##           Model Test.Accuracy
## 1 Random Forest          51.16
## 2 Naive Bayes            42.80
## 3 KNN                    41.04
```

## Association rule mining

In this section, I will use the data on grocery purchases in `groceries.txt` and find some interesting association rules for these shopping baskets.

```
# plot the most frequent items (top 10)
itemFrequencyPlot(tdata, topN = 10)
```



```
# Visualizing the rules
inspect(sort(associa_rules, by = 'lift')[1:10])
```

```
##      lhs      rhs      support      confidence coverage      lift
## [1] {salty}    => {snack}    0.003050330 0.7692308  0.003965430 51.46520
## [2] {misc.}    => {beverages} 0.003152008 0.6326531  0.004982206 51.42267
## [3] {snack,long} => {bakery}    0.003355363 1.0000000  0.003355363 26.72554
## [4] {snack,long} => {life}      0.003355363 1.0000000  0.003355363 26.72554
## [5] {juice,long} => {bakery}    0.004270463 1.0000000  0.004270463 26.72554
## [6] {juice,long} => {life}      0.004270463 1.0000000  0.004270463 26.72554
## [7] {product}  => {bakery}    0.013218099 1.0000000  0.013218099 26.72554
## [8] {bakery}    => {product}    0.013218099 0.3532609  0.037417387 26.72554
## [9] {product}  => {life}      0.013218099 1.0000000  0.013218099 26.72554
## [10] {life}     => {product}    0.013218099 0.3532609  0.037417387 26.72554
##      count
## [1] 30
## [2] 31
## [3] 33
## [4] 33
```

```
plot(associa_rules, method = "graph",
     measure = "confidence", shading = "lift")
```

size: confidence (0.3 – 1)  
color: lift (3.139 – 51.465)

