

Homework, The Last

Aaron Politsky

December 3, 2015

Part 1

After we load up the emails and create a Corpus, let's preprocess it and create the document term matrix.

```
corpus = tm_map(corpus, content_transformer(removeNumbers))
corpus = tm_map(corpus, content_transformer(removePunctuation))
corpus = tm_map(corpus, content_transformer(tolower))
corpus = tm_map(corpus, content_transformer(removeWords), stopwords("english"))
corpus = tm_map(corpus, content_transformer(stripWhitespace))

dtm = DocumentTermMatrix(corpus)
```

We'll use a sparsity value of .99 at first

```
# next we remove infrequent words
# we keep columns that are less than 0.99 percent sparse
# this is the parameter that you will need to tune in the homework
sparse_dtm = removeSparseTerms(x=dtm, sparse = 0.99)
sparse_dtm
```

```
## <<DocumentTermMatrix (documents: 5728, terms: 1636)>>
## Non-/sparse entries: 331057/9039951
## Sparsity          : 96%
## Maximal term length: 16
## Weighting          : term frequency (tf)
```

```
# convert all elements to binary
# The occurrence of the word fantastic tells us a lot
# The fact that it occurs 5 times may not tell us much more
sparse_dtm = weightBin(sparse_dtm)
```

```
# split into train and test using sampling_vector
df = as.data.frame(as.matrix(sparse_dtm))
df_train = df[sampling_vector,]
df_test = df[-sampling_vector,]
spam_train = emails$spam[sampling_vector]
spam_test = emails$spam[-sampling_vector]
```

```
library(e1071)
```

Naive Bayes

Let's train our classification model:

```

### your code for classification goes below
nb_model = naiveBayes(df_train, spam_train)

if (file.exists("nb_train_predictions.RData")) {
  load("nb_train_predictions.RData")
} else {
  nb_train_predictions = predict(nb_model, df_train)
  save(nb_train_predictions, file = "nb_train_predictions.RData")
}
mean(nb_train_predictions == spam_train)

```

```
## [1] 0.9024443
```

```
table(actual = spam_train, predictions = nb_train_predictions)
```

```

##      predictions
## actual  ham spam
##   ham 3071  431
##   spam   16 1064

```

```

# compute test error
if (file.exists("nb_test_predictions.RData")) {
  load("nb_test_predictions.RData")
} else {
  nb_test_predictions = predict(nb_model, df_test)
  save(nb_test_predictions, file = "nb_test_predictions.RData")
}

```

Our test accuracy is:

```
mean(nb_test_predictions == spam_test)
```

```
## [1] 0.8970332
```

```
table(actual = spam_test, predictions = nb_test_predictions)
```

```

##      predictions
## actual  ham spam
##   ham  747  111
##   spam   7  281

```

Stemming

Let's try Stemming

```

# also try stemming as a preprocessing step
stemmed = tm_map(corpus, stemDocument, language = "english")
stemmed.dtm = DocumentTermMatrix(stemmed)
stemmed.dtm = removeSparseTerms(x=stemmed.dtm, sparse = 0.99)

```

```

stemmed.dtm = weightBin(stemmed.dtm)
stemmed_df = as.data.frame(as.matrix(stemmed.dtm))
stemmed_df_train = stemmed_df[sampling_vector,]
stemmed_df_test = stemmed_df[-sampling_vector,]

# train model on stemmed corpora
model_stem = naiveBayes(stemmed_df_train, spam_train)
if (file.exists("nb_test_predictions_stem.RData")) {
  load("nb_test_predictions_stem.RData")
} else {
  nb_test_predictions_stem = predict(model_stem, stemmed_df_test)
  save(nb_test_predictions_stem, file = "nb_test_predictions_stem.RData")
}

```

Our test accuracy using stemming is:

```
mean(nb_test_predictions_stem == spam_test)
```

```
## [1] 0.9144852
```

```
table(actual = spam_test, predictions = nb_test_predictions_stem)
```

```
##      predictions
## actual ham spam
##   ham 769   89
##   spam   9 279
```

Different values of Sparsity

Let's try a range of sparsity values:

```

# try different values of sparsity
stemmed.dtm = DocumentTermMatrix(stemmed)
sparsity <- seq(0.9, 0.99, by = .02)
dtms.by.sparsity <- lapply(sparsity, function(sp) {
  as.data.frame(as.matrix(weightBin(removeSparseTerms(x=stemmed.dtm, sparse = sp))))
})

if (file.exists("nb.test.predictions.by.sparsity.Rda")) {
  load("nb.test.predictions.by.sparsity.Rda")
} else {
  nb.test.predictions.by.sparsity <-
    lapply(dtms.by.sparsity, function(dtm) {
      df.train = dtm[sampling_vector,]
      df.test  = dtm[-sampling_vector,]
      nb_model = naiveBayes(df.train, spam_train)
      nb_test_predictions = predict(nb_model, df.test)
    })
  save(nb.test.predictions.by.sparsity, file = "nb.test.predictions.by.sparsity.Rda")
}
names(nb.test.predictions.by.sparsity) <- sparsity

```

Our test accuracy as a function of sparsity is:

```
sapply(nb.test.predictions.by.sparsity, function(nb_test_predictions)
  mean(nb_test_predictions == spam_test)
)
```

```
##      0.9      0.92      0.94      0.96      0.98
## 0.8952880 0.8926702 0.8952880 0.9162304 0.9354276
```

Part 2

Specify a stochastic block matrix: We believe fraudsters don't often trade with themselves or with honest, but mostly with accomplices.

Accomplices try to look normal by trading a lot with honest people, and occasionally with fraudsters. Honest trade with accomplices and themselves, but rarely with fraudsters.

```
# stochastic block matrix
# our graph will have three communities
eps <- .021 # epsilon

M = matrix(c(
  eps,      1-2*eps,      eps,
  1-2*eps,   2*eps,   .5-2*eps,
  eps,      .5-2*eps,   .5-2*eps
),
nrow = 3)
M
```

```
##      [,1] [,2] [,3]
## [1,] 0.021 0.958 0.021
## [2,] 0.958 0.042 0.458
## [3,] 0.021 0.458 0.458
```

These are symmetric:

```
t(M) == M
```

```
##      [,1] [,2] [,3]
## [1,] TRUE TRUE TRUE
## [2,] TRUE TRUE TRUE
## [3,] TRUE TRUE TRUE
```

Let's build our graph

```
# sample a random graph
# 100 nodes grouped into 3 communities
num.fraudsters <- num.accomplices <- 5
num.honest <- 90
rg =
  sample_sbm(num.fraudsters + num.accomplices + num.honest, # number of nodes
```

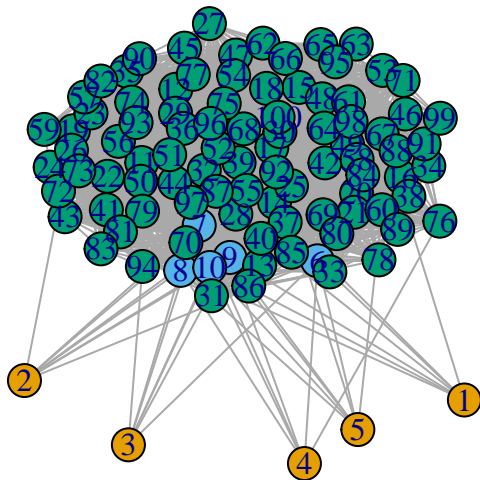
```

    pref.matrix = M,
    block.sizes = c(num.fraudsters,
                    num.accomplices,
                    num.honest), # how many nodes belong to each community
    loops = F,                  # no loops (vertex that connects to itself)
    directed = F                # we want an undirected graph
  )

# membership vector used to color vertices
membership_vector = c(rep(1, num.fraudsters),
                      rep(2, num.accomplices),
                      rep(3, num.honest))

plot_layout = layout.fruchterman.reingold(rg)
plot(rg,
     vertex.color=membership_vector,
     layout = plot_layout)

```



Our plot shows that the fraudsters are separate from the rest, but it's hard to identify the honest from accomplices, which is the point of accomplices.

```

library(lda)

result =
  mmsb.collapsed.gibbs.sampler(get.adjacency(rg), # the adjacency matrix
                               K = 3,             # 3 groups
                               num.iterations=10000,
                               alpha = 0.1,
                               burnin = 500L,
                               beta.prior = list(1, 1))

# this matrix tells us for each vertex what is the probability that it belongs to
# one of the K communities
memberships = with(result, t(document_sums) / colSums(document_sums))
head(memberships,20)

##           [,1]      [,2]      [,3]
## [1,] 0.000000000 0.00000000 1.000000000
## [2,] 0.000000000 0.00000000 1.000000000

```

```
## [3,] 0.000000000 0.00000000 1.000000000
## [4,] 0.000000000 0.00000000 1.000000000
## [5,] 0.015151515 0.00000000 0.984848485
## [6,] 0.969696970 0.00000000 0.030303030
## [7,] 0.969696970 0.03030303 0.000000000
## [8,] 1.000000000 0.00000000 0.000000000
## [9,] 0.989898990 0.01010101 0.000000000
## [10,] 1.000000000 0.00000000 0.000000000
## [11,] 0.015151515 0.98484848 0.000000000
## [12,] 0.000000000 0.99494949 0.005050505
## [13,] 0.000000000 0.80808081 0.191919192
## [14,] 0.000000000 1.00000000 0.000000000
## [15,] 0.075757576 0.92424242 0.000000000
## [16,] 0.000000000 1.00000000 0.000000000
## [17,] 0.000000000 1.00000000 0.000000000
## [18,] 0.005050505 0.99494949 0.000000000
## [19,] 0.429292929 0.57070707 0.000000000
## [20,] 0.000000000 0.99494949 0.005050505
```

Many of these are clearly predicted to be from one community versus the others, indicating that we think they likely belong to a given community.

We expect the first five to be fraudsters, group 1:

```
head(memberships[1:5,])
```

```
##           [,1] [,2]      [,3]
## [1,] 0.00000000    0 1.0000000
## [2,] 0.00000000    0 1.0000000
## [3,] 0.00000000    0 1.0000000
## [4,] 0.00000000    0 1.0000000
## [5,] 0.01515152    0 0.9848485
```

We see that they are classified together.

And the expected accomplices:

```
memberships[6:10,]
```

```
##           [,1]      [,2]      [,3]
## [1,] 0.969697 0.00000000 0.03030303
## [2,] 0.969697 0.03030303 0.00000000
## [3,] 1.000000 0.00000000 0.00000000
## [4,] 0.989899 0.01010101 0.00000000
## [5,] 1.000000 0.00000000 0.00000000
```

These are mainly classified together as well.

The honest ones?

```
memberships[11:30,]
```

```
##           [,1]      [,2]      [,3]
## [1,] 0.015151515 0.9848485 0.000000000
## [2,] 0.000000000 0.9949495 0.005050505
## [3,] 0.000000000 0.8080808 0.191919192
## [4,] 0.000000000 1.0000000 0.000000000
## [5,] 0.075757576 0.9242424 0.000000000
## [6,] 0.000000000 1.0000000 0.000000000
## [7,] 0.000000000 1.0000000 0.000000000
## [8,] 0.005050505 0.9949495 0.000000000
## [9,] 0.429292929 0.5707071 0.000000000
## [10,] 0.000000000 0.9949495 0.005050505
## [11,] 0.000000000 0.9292929 0.070707071
## [12,] 0.000000000 1.0000000 0.000000000
## [13,] 0.000000000 0.9494949 0.050505051
## [14,] 0.343434343 0.5808081 0.075757576
## [15,] 0.212121212 0.7878788 0.000000000
## [16,] 0.005050505 0.8434343 0.151515152
## [17,] 0.000000000 1.0000000 0.000000000
## [18,] 0.227272727 0.7727273 0.000000000
## [19,] 0.000000000 1.0000000 0.000000000
## [20,] 0.035353535 0.9646465 0.000000000
```

```
mean(memberships[11:100,1])
```

```
## [1] 0.04393939
```

```
mean(memberships[11:100,2])
```

```
## [1] 0.9083053
```

```
mean(memberships[11:100,3])
```

```
## [1] 0.04775533
```

The honest ones look to be grouped as well.

```
# the estimate of the stochastic block matrix M
ratio = with(result, blocks.pos / (blocks.pos + blocks.neg))
ratio
```

```
##           [,1]      [,2]      [,3]
## [1,] 0.0250000 0.475476839 0.931506849
## [2,] 0.5139276 0.492738275 0.002594034
## [3,] 0.7792208 0.001331558 0.000000000
```

```
# actual M
M
```

```
##           [,1] [,2] [,3]
## [1,] 0.021 0.958 0.021
## [2,] 0.958 0.042 0.458
## [3,] 0.021 0.458 0.458
```

Community assignments got shifted around, but we classified consistently. If we were to rearrange one of these matrices so the groupings were consistent, we might see them line up well.