

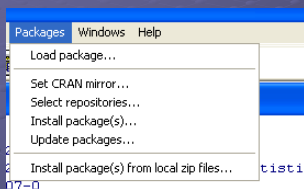
Using R for Text Mining



Part 2 : R Text Mining

Assumes knowledge of R but not how to text mine in R

Install the tm package



TM – source documents

```
#####  
> # read in source document collection  
#####  
> txt.csv <- read.csv(file="c:/text mining/Top2lss.txt", header=FALSE)  
> txt <- Corpus(DataframeSource(txt.csv))  
  
> summary(txt)  
A corpus with 330 text documents  
  
# examine the first 10 rows  
> inspect(txt[1:10])  
[[1]]  
A crisis that could affect our ability to regulate ourselves.  
  
[[2]]  
A need to deal more thoroughly with non-traditional risk management approaches  
  
[[3]]  
Ability of members to prove they are more than just number crunchers  
  
[[4]]  
ability to convince non-insurance companies of the value/skills offered by CAS members.
```

TM – preprocess lowercase

```
> #####  
> # a little pre-processing to prep the data for TM  
> # convert to lower case  
> # tmTolower is one of several available text transformations.  
> # To see all currently available use: getTransformations()  
> #####  
> txt <- tm_map(txt, tolower)  
  
> inspect(txt[1:10])  
[[1]]  
a crisis that could affect our ability to regulate ourselves.  
  
[[2]]  
a need to deal more thoroughly with non-traditional risk management approaches  
  
[[3]]  
ability of members to prove they are more than just number crunchers  
  
[[4]]  
ability to convince non-insurance companies of the value/skills offered by cas members.
```

TM – search & replace

```
#####  
> # Replace the slashes in the text with a blank  
> #  
#####  
> # txt <- gsub("/", " ", txt)  
> for (j in 1:length(txt)) txt[[j]] <- gsub("/", " ", txt[[j]])  
  
> inspect(txt[1:10])  
[[1]]  
a crisis that could affect our ability to regulate ourselves.  
  
[[2]]  
a need to deal more thoroughly with non-traditional risk management approaches  
  
[[3]]  
ability of members to prove they are more than just number crunchers  
  
[[4]]  
ability to convince non-insurance companies of the value skills offered by cas members.
```

```

TM – search & replace
con't
> # Replace other characters
> #
> for (j in 1:length(txt)) txt[[j]] <- gsub("[&-\^()\\]", " ", txt[[j]])
> for (j in 1:length(txt)) txt[[j]] <- gsub("[&|+|\\(|\\)|\\.|,|' ", txt[[j]]);
> inspect(txt[1:10])
[[1]]
a crisis that could affect our ability to regulate ourselves

[[2]]
a need to deal more thoroughly with non traditional risk management approaches

[[3]]
ability of members to prove they are more than just number crunchers

[[4]]
ability to convince non insurance companies of the value skills offered by cas members

[[5]]
ability to help sort out property pricing problems

```

```

TM – search & replace con't
> # Replace enterprise risk management
> #
> for (j in 1:length(txt)) txt[[j]] <- gsub("enterprise risk management", "erm", txt[[j]])
> for (j in 1:length(txt)) txt[[j]] <- gsub("off shoring", "offshoring", txt[[j]]);
> inspect(txt[1:10])
[[1]]
a crisis that could affect our ability to regulate ourselves

[[2]]
a need to deal more thoroughly with non traditional risk management approaches

[[3]]
ability of members to prove they are more than just number crunchers

[[4]]
ability to convince non insurance companies of the value skills offered by cas members

[[5]]
ability to help sort out property pricing problems

```

```

TM – search & replace con't
> # remove stopwords
> #
> txt <- tm_map(txt, removeWords, stopwords("english"))
> #ublnit(dbe=txtcsv)
>
> # remove punctuation
> #
> txt <- tm_map(txt, removeNumbers)
> txt <- tm_map(txt, removePunctuation)
> inspect(txt[1:10])
[[1]]
crisis affect ability regulate

[[2]]
deal thoroughly traditional risk management approaches

[[3]]
ability prove crunchers

[[4]]
ability convince insurance companies value skills offered cas

[[5]]
ability help sort property pricing

```

TM – search & replace con't

```
> # remove stopwords & punctuation
> #
> txt <- tm_map(txt, removeWords, stopwords("english"))
> txt <- tm_map(txt, removeNumbers)
> txt <- tm_map(txt, removePunctuation)
> inspect(txt[1:10])
[[1]]
crisis affect ability regulate

[[2]]
deal thoroughly traditional risk management approaches

[[3]]
ability prove crunchers

[[4]]
ability convince insurance companies value skills offered cas

[[5]]
ability help sort property pricing
> length(txt)
[1] 330
```

TM – search & replace con't

```
(j in 1:length(txt)) txt[[j]] <- gsub("professional", "professions", txt[[j]]);
>
> getTransformations()
[1] "as.PlainTextDocument" "convert_UTF_8" "removeNumbers"
"removePunctuation"
[5] "removeWords" "stemDocument" "stripWhitespace"
> txt <- tm_map(txt, stemDocument)
> inspect(txt[40:50])
[[1]]
climat chang

[[2]]
compani complain oner expen educ system set cas

[[3]]
compani will pay meet

[[4]]
compet profess organ

[[5]]
competit profess
```

TM – Document by Term Matrix

```
> # create a document by term matrix
> #
> dtm <- DocumentTermMatrix(txt)
> nrow(dtm), ncol(dtm)
[1] 330
[1] 469
> inspect(dtm[1:24,1:9])
A document-term matrix (24 documents, 9 terms)

Non-sparsity entries: 10/206
Sparsity : 95%
Maximal term length: 8
Weighting : term frequency (tf)

 abil aca accept account accredit accuraci activ actuari address
1 1 0 0 0 0 0 0 0 0 0
2 0 0 0 0 0 0 0 0 0
3 1 0 0 0 0 0 0 0 0
4 1 0 0 0 0 0 0 0 0
5 1 0 0 0 0 0 0 0 0
6 0 0 1 0 0 0 0 0 0
7 0 0 0 0 0 0 0 1 0
8 0 0 0 0 0 0 0 1 0
9 0 0 0 0 0 0 0 1 0
10 0 0 0 0 0 0 0 0 0
```

Three Variations on Tag Clouds

```
#####
# install.packages("fun", repos="http://r-forge.r-project.org").
# with style
#####
require(fun)
data(tagData)
v <- as.matrix(sort(sapply(top2,
doit),decreasing=TRUE)[1:numwords],
colnames=count)v[,1:numwords]v
x <- data.frame(rownames(v), "http://www.casact.org",
tagData$color[1:length(v)],tagData$icolor[1:length(v)],
colnames(x) <- c(
'count','color','icolor');x
htmlFile=paste(tempfile(), ".html", sep="")
#htmlFile=paste("tagData", ".html", sep="")
if (file.create(htmlFile)) {
tagCloud(x, htmlFile)
browseURL(htmlFile)
}
```



Get Individual Records associated with the term "actuar"

```
> #####
> # Get records for the word actuar
> #####
> subset(top2, actuar=1)[1:10,1:10]
  abil account actuari cas casualti chang compani compett continu credibil
7  0  0  1  0  0  0  0  0  0  0  0
8  0  0  1  0  0  0  0  0  0  0  0
9  0  0  1  0  0  0  0  0  0  0  0
17 0  0  1  0  0  0  0  0  0  0  0
> rownames(subset(top2, actuar=1)) #which rows have actuar=1
[1] "7" "8" "9" "17" "20" "45" "46" "57" "59" "71" "74" "77" "78" "81"
[6] "82" "95" "105" "118" "136" "149" "153" "160" "163" "173" "180" "182" "183" "186"
[20] "191" "193" "195" "196" "197" "204" "205" "206" "214" "220" "221" "223" "224" "225"
[43] "227" "229" "250" "261" "263" "266" "277" "282" "298" "299" "309" "313" "318" "321"
[57] "323" "328" "329"
> txt.actuari <- txt.csv(rownames(subset(top2, actuar=1)),)
> txt.actuari[1:5], #print out records with actuar
[1] Actuarial Malpractice (2)
[2] Actuarial Students of today are not good communicators/executive material.
[3] Actuaries regain status as masters of risk
[4] automation/modeling of actuarial skills
[5] Being a market leading brand of actuarial sciences - what do we stand for?
330 levels: approaches to address external changes ERM ... Will we get to a point where we have too
many members??
```

Correlation Plot

```
> #####
> # produce a Correlation plot with the most frequently occurring top 24 terms
> #####
> require(ellipse)
> nrow(dtm3); ncol(dtm3)
[1] 330
[1] 36
> dtm4 <- removeSparseTerms(dtm, 0.97)
> nrow(dtm4); ncol(dtm4)
[1] 330
[1] 23
> inspect(dtm4[1:3,])

actuari cas chang compani compett credibil current educ erm exam financ increa industri
1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
2  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
3  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
insur intern manag model organ profess regul research reserv standard
1  0  0  0  0  0  0  0  1  0  0  0  0
2  0  0  1  0  0  0  0  0  0  0  0
3  0  0  0  0  0  0  0  0  0  0  0  0
> top4.matrix <- as.matrix(dtm4)
> txt.cor <- cor(top4.matrix)
> ord <- order(txt.cor[lower.tri(dtm4)])
> xc <- txt.cor[ord, ord]
> plotcorr(xc, col=cm.colors(12)[5*xc + 6])
>
```

Principal Components

- An unsupervised technique
- Groups similar variables (rather than similar records) together
- Uses correlation matrix – variables (here terms) that are highly correlated are grouped together

R Princomp function

```
>require(stats)
>prcomp(top2, scale=TRUE)
>summary(prcomp(top2, scale=TRUE))
● Then examine top components
>.for (i in 1:15) {
> top4[[i]] <- sort(survey.prcomp$rotation[,i],
decreasing=TRUE)[1:4]}
>top4.
```

The Top Components

Component	Concept
1	practical research
2	predictive modeling
3	maintain standards
4	demand for actuarial skills
5	risk management
6	exam structure
7	actuarial organization
8	reputation/financial issues
9	economic capital
10	competition - other professions
11	opportunities for actuaries
12	capital
13	predictive modeling
14	globalism/competition
15	reserving and credibility
