

# Open Source Text Mining

Mathew Flynn, PhD  
Louise Francis, FCAS, MAAA

---

---

---

---

---

---

---

---

## Rationale For Paper

- Text mining is a new and promising technology for analyzing unstructured text data
- Commercial text mining software can be expensive and difficult to learn
- Several free open source languages can perform text mining, but without help they can be difficult to learn
- In this session we will provide a mini- tutorial to 2 open source products

---

---

---

---

---

---

---

---

## Two Open Source Products for Text Mining

- Perl – a text processing language
- R – a statistical and analytical language with text mining functionality provided by a text mining package tm

---

---

---

---

---

---

---

---

## The Data

- Text mining can be applied to many common tasks
  - Internet searches
  - Screening emails for spam
  - Analyzing free form fields in underwriting and claims files
  - Analyzing survey data
- We illustrate the last 2
- Survey data can be downloaded from CAS web site.

---

---

---

---

---

---

---

---

## Mini Tutorial

- We will give tutorial on using Perl and R for text mining
- Download the survey data
- Follow our examples

---

---

---

---

---

---

---

---

## The Survey Data

- From 2008 CAS Quinquennial Survey
- What are the top two issues that will impact the CAS in the next five years?

Survey Question: Top Two Issues Affecting CAS

---

A crisis that could affect our ability to "regulate" ourselves.

---

A need to deal more thoroughly with non-traditional risk management approaches

---

Ability of members to prove they are more than just number crunchers

---

ability to convince non-insurance companies of the value/skills offered by CAS members.

---

---

---

---

---

---

---

---

---

## Perl

- Go to [www.Perl.org](http://www.Perl.org)
- Download Perl
- Run execute file (or run active perl)
- the Windows command search must be correct for Windows to find the desired perl.exe

---

---

---

---

---

---

---

---

## www.Perl.org



---

---

---

---

---

---

---

---

## Good References

- *Practical Text Mining with Perl* by Bilisoly (2008) is an excellent resource for text mining in Perl
- *Perl for Dummies* (Hoffman, 2003) provides a basic introduction including needed header information

---

---

---

---

---

---

---

---

## Some Key Things

- Perl must be run from DOS. One gets to DOS by finding the Command Prompt on the Programs menu
- Before running Perl switch to the Perl directory (i.e., if Perl was installed and it is in the folder named Perl, in DOS, type "cd C:\Perl").
- Programs need to be saved in text processing software. We recommend Notepad rather than Word, as some of the features of Word cause unexpected results when running a program. We recommend using the extension .pl.

---

---

---

---

---

---

---

---

## Some Key Things cont.

- The header line of a Perl program is dependent on the operating system.
- To run a Perl program type the following at the command prompt:
  - Perl program\_name input\_file\_name, output\_file\_name[1]
- The input and output files are only required if the program requires a file, and the file name is not contained in the program itself.
- The input and output file may be contained in the program code, as illustrated in some of our examples, rather than entered at runtime.

---

---

---

---

---

---

---

---

## Parsing Text

- Identify the spaces, punctuation and other non alphanumeric characters found in text documents and separating the words from these other characters
- Most computer languages (and spreadsheets) have text functions that perform the search and substring functions to do this
- Perl has special functions for parsing text

---

---

---

---

---

---

---

---

## The split function

- `split(/separating character(s)/, string)`
- Example
  - `$Response = "Ability of members to prove they are more than just number crunchers";`
  - `@words =split (/ /, $Response);`

---

---

---

---

---

---

---

---

## Complications of split function

- More than one space
  - `@words =split (/ [s+]/, $Response);`
- Other separators
  - Use substitute function

---

---

---

---

---

---

---

---

## Simple parse program: Parse2.pl

- `#!/perl -w`
- `# Parse2.pl`
- `# Program to parse text string using one or more spaces as separator`
- `$Response = "Ability of members to prove they are more than just number crunchers";`
- `@words =split (/s+/, $Response); #parse words in string`
- `# Loop through words in word array and print them`
- `foreach $word (@words) {`
- `print "$word\n";`
- `}`

---

---

---

---

---

---

---

---

## Less Simple parse program: Parse3.pl

```
#!/perl -w
# Parse3.pl
# Program to parse a sentence and remove punctuation
$Test = "A crisis that could affect our ability to 'regulate'
ourselves.;"# a test string with punctuation
@words =split (/(\s+)/, $Test); # parse the string using spaces
# Loop through words to find non punctuation characters
foreach $word (@words) {
  while ( $word =~ /(\w+)/g ) {
    # match by 1 or more alphanumeric characters. These will be the
    words excluding punctuation
    print "$1 \n"; #print the first match which will be the word of
    alphanumeric characters
  }
}
```

---

---

---

---

---

---

---

---

## Read in survey data and parse

```
#!/perl -w
# Enter file name withtext data here
$TheFile = "Top2lss.txt";
# open the file
open(INFILE, $TheFile) or die "File not found";
# read in one line at a time
while(<INFILE>) {
  chomp; # eliminate end of line character
  s/[?!'(){}&|]/g; # replace punctuation with null
  s// /g; # replace slash with space
  s/\-/ /g; #replace dash with null
  s/^\s/g; #replace beginning of line space
  print "$ \n"; # print cleaned line out
  @word=split(/(\s+)/); # parse line
  _
}
```

---

---

---

---

---

---

---

---

## Print it out also

```
#!/perl -w
# parsecomplex.pl
# Enter file name withtext data here
$TheFile = "Top2lss.txt";
# open the file
open(INFILE, $TheFile) or die "File not found";
# read in one line at a time
while(<INFILE>) {
  chomp; # eliminate end of line character
  s/[?!'(){}&|]/g; # replace punctuation with null
  s// /g; # replace slash with space
  s/\-/ /g; #replace dash with null
  s/^\s/g; #replace beginning of line space
  print "$ \n"; # print cleaned line out
  @word=split(/(\s+)/); # parse line
  foreach $word (@word) {
    print "$word,";
    print "\n";
  }
}
```

---

---

---

---

---

---

---

---

## Word Search

- First, read in the accident description field
- For each claim
  - Read in each word
  - If the lower case of the target word is found output a 1 for the new indicator variable, otherwise output a 0.

---

---

---

---

---

---

---

---

## SearchTarget.pl

```
SearchTarget.pl
# Open an input stream
# Substitute the variables containing the web host data
# Write the word list to a file
# Read the file and store the words in an array
# Print the words to the screen
# Read the accident description field
# Read each word
# Check if the word is in the array
# If found, set the indicator variable to 1
# If not found, set the indicator variable to 0
# Print the indicator variable to the screen
# Close the input stream
```

---

---

---

---

---

---

---

---

## Using Target in Analysis

Homeowner Claim	Mean Severity
No	2,376.6
Yes	6,221.1

---

---

---

---

---

---

---

---

## Text Statistics

- The length of each word is tabulated within a loop. A key line of code is:
- `$count[length($x)] +=1; #increment counter for words of this length`

---

---

---

---

---

---

---

---

## Perl Program for Word Lengths

```
#!/usr/bin/perl
# Error file name with text data here
my $infile = "top500.txt";
# open file
open(INFILE, $infile) or die "File not found";
# read in one line at a time
while(<INFILE>){
    # change # terminate end of line character
    s/[\n]//g; # replace punctuation with null
    s/ / /g; # replace slash with space
    s/"/ /g; # replace double quote with null
    s/"/ /g; # replace beginning of line space
    print "in" # print cleaned line out
    @words = split(/\s+/); # parse line
    # count length of each word in array @count
    foreach $x (@words){
        $count{length($x)} +=1;
    }
    $mcount = $count;
    # print out largest word size and frequency of each count
    print "Count $mcount\n";
    for ($i = 0; $i <= $count;){
        # show word of that size exist
        if (exists($count{$i})) {
            print "There are $count{$i} words of length $i\n";
        }
        $i += 1; # increment loop counter
    }
}
```

---

---

---

---

---

---

---

---

## Hashes

- A hash is like an array, but can be distinguished from an array in a number of ways.
- An array is typically indexed with zero and integer values, while a hash can be indexed with a letter or word.
- Instead of an index the hash has a key that maps to a specific array value.
  - For instance, while the first entry in a Perl array is `$array[0]`, the first entry in a hash might be `$hash{'a'}` or `$hash{'x'}` or even `$hash{'hello'}`. (Note that the order is not relevant.)
- Because the time to locate a value on a hash table is independent of its size, hashes can be very efficient for processing large amounts of data

---

---

---

---

---

---

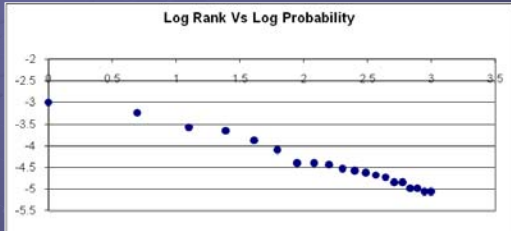
---

---





## Zipf's Law



---

---

---

---

---

---

---

---

## Stop Words

- Frequently occurring words
- The
- A
- To
- It
- Do not contribute to meaning of record of text
- Eliminate

---

---

---

---

---

---

---

---

## Substitution operator

- Thus to eliminate the word "the", use the code

`s/the//g;`

- Apply to multiple terms you want to eliminate

`s/[-?!'(){}&.;]/g;`

---

---

---

---

---

---

---

---

## Term Document Matrix

- A Table of indicator variables
- If a word is present, a 1, otherwise a 0

---

---

---

---

---

---

---

---

## Term Data Matrix

	Ourselves	cas	Not	That	communicators	executive	our	approaches
1	0	0	1			0	1	0
0	0	0	0			0	0	1
0	0	0	0			0	0	0
0	1	0	0			0	0	0
0	0	0	0			0	0	0
0	0	0	0			0	0	0
0	0	0	0			0	0	0
0	0	1	0			1	0	0
0	0	0	0			0	0	0
0	0	0	0			0	0	0

---

---

---

---

---

---

---

---

## Word Lengths

Length	GL	
	Data	Survey Data
1	1,062	21
2	4,172	309
3	5,258	298
4	5,418	215
5	2,982	153
6	2,312	143
7	2,833	213
8	1,572	161
9	1,048	216
10	591	146
11	111	92
12	156	44
13	78	61
14	19	2
15	0	3
16	1	1
17	2	0
18	1	0
19	1	0

---

---

---

---

---

---

---

---

