# Assignment Ten

## Sybil Melton

December 10, 2014

# CONTENTS

# LIST OF TABLES

# LISTINGS

# 1 BLOG RETRIEVAL

Choose a blog or a newsfeed (or something similar as long as it has an Atom or RSS feed). It should be on a topic or topics of which you are qualified to provide classification training data. In other words, choose something that you enjoy and are knowledgable of. Find a feed with at least 100 entries. Create between four and eight different categories for the entries in the feed.
Download and process the pages of the feed as per the week 12 class slides.

## 1.1 SOLUTION

To find a blog, I used the "next blog" functionality to randomly choose one until I found one with a technology topic and more than 100 entries. The particular blog I chose is titled "The Invisible Things Lab's blog." The description is Kernel, Hypervisor, Virtualization, Trusted Computing and other system-level security stuff, which I have knowledge of from work and my studies at ODU. I used requests, feedparser, and BeautifulSoup to get the blog pages and parse the feed. I retrieved the first four pages and saved the entries in a list, to be filtered later. The content still had some markup language and illegal XML characters, so I used the same functions I created for assignment nine to remove them. Listing 1 is the code used to accomplish the first part of this task. [4]

```
 1  entries = []
 2  blog = 'http://theinvisiblethings.blogspot.com/'
 3  path = '/feeds/posts/default'
 4  query_arg = {'alt' : 'atom'}
 5  udata = urllib.urlencode(query_arg)
 6
 7  full_url = blog+path+"?"+udata
 8  r = requests.get(full_url)
 9  f=feedparser.parse(full_url)
10  entries = f['entries']
11  bsdata = r.text
12  soup = BeautifulSoup(bsdata)
13  next = soup.find('link', rel='next')
14  for x in range(1,4):
15    n= next.get('href')
16    f=feedparser.parse(n)
17    entries.extend(f['entries'])
18    r = requests.get(n)
19    bsdata = r.text
20    soup = BeautifulSoup(bsdata)
21    next = soup.find('link', rel='next')
22
23  _illegal_xml_chars_RE = re.compile(u'[\x00-\x08\x0b\xa9\xae\x0c\x0e-\x1F\uD800-\uDFFF\uFFFE\uFFFF]')
24
25  def remove_tags(text):
26      cleanr =re.compile('<.*?>')
27      cleantext = re.sub(cleanr,'', text)
28      cleantext = re.sub("\n", ' ', cleantext)
29      cleantext = cleantext.strip()
30      cleantext = cleantext.replace(u'\u201c', '')
31      cleantext = cleantext.replace(u'\u201d', '')
32      cleantext = cleantext.replace(u'\u2019', '\'')
33      cleantext = cleantext.replace('-&gt;', '')
34      return _illegal_xml_chars_RE.sub(' ', cleantext)
35
36  for x in range(0,100):
37    entries[x]['content'][0]['value'] = remove_tags(entries[x]['content'][0]['value'])
38    entries[x]['title'] = remove_tags(entries[x]['title'])
```

Listing 1: Retrieve Blog

As I read the entries, I chose announcements, security, hardware, os, and paper as the categories. Announcements was for general announcements. Security was for any entry about cybersecurity. Hardware entries discussed chipset and DRAM. Os handled any general articles about operating systems that did not fit into security, and paper handled any blogs that just posted white papers or slides from conferences.

## 2 CLASSIFICATION

Manually classify the first 50 entries, and then classify (using the fisher classifier) the remaining 50 entries. Report the cprob() values for the 50 titles as well. From the title or entry itself, specify the 1-, 2-, or 3-gram that you used for the string to classify. Do not repeat strings; you will have 50 unique strings.

Create a table with the title, the string used for classification, cprob(), predicted category, and actual category.

### 2.1 SOLUTION

My program set up the fisher classifier and database. The database was "feed.db" and is included with this report. Then feedfilter was used to manually classify the first 50 entries, train the classifier, and guess a classification for the remaining 50. Listing 2 is the code used to accomplish the next task. [4]

```
1  cl=docclass.fisherclassifier(docclass.getwords)
2  cl.setdb('feed.db')
3
4  feedfilter.read(entries, cl)
```

Listing 2: Set Up Classifier

I modified feedfilter.py to classify the entries. I changed the read function to accept my list of entries, so the list could be split between the first and second 50 entries. The 'summary' token was changed to the content value, because I found some entries were cut off because of the length. The title and content were concatenated before sent to the classifier. As the loop progressed, it added my classifier and manually assigned categories to the list for each entry. Then the train function from docclass.py was called and was not modified. Listing 3 is the modified code in feedfilter.py to accomplish this. [4]

```
1   for x in range(0, 50):
2     print
3     print '-----'
4     # Print the contents of the entry
5     print 'Title:    '+entries[x]['title'].encode('utf-8')
6     print
7     print entries[x]['content'][0]['value'].encode('utf-8')
8
9     # Combine all the text to create one item for the classifier
10    fulltext= entries[x]['title'] + ', ' + entries[x]['content'][0]['value']
11
12    # Ask the user to specify the correct category and train on that
13    tr=raw_input('Enter classifier: ')
14    entries[x]['classifier'] = tr
15    cl=raw_input('Enter category: ')
16    entries[x]['actual'] = cl
17    classifier.train(fulltext,cl)
```

Listing 3: First 50 Classification

Next the second 50 entries had to classified by me and the classifier. Additionally, the cprob had to be saved, so I modified docclass.py classify function to return it when a guess was given. This allowed me to save both the predicted category and cprob for each in the entries list. Listing 4 is the additional code in feedfilter.py to accomplish this. [4]

```
1  # classify the final 50 and get the guess and cprob
2    for x in range(50, 100):
3    print
4    print '-----'
5    # Print the contents of the entry
6    print 'Title:    '+entries[x]['title'].encode('utf-8')
7    print
8    print entries[x]['content'][0]['value'].encode('utf-8')
9
```

```
10    # Combine all the text to create one item for the classifier
11    fulltext= entries[x]['title'] + ', ' + entries[x]['content'][0]['value']
12
13    # Ask the user to specify the correct category
14    cl=raw_input('Enter category: ')
15    entries[x]['actual'] = cl
16    # Print the best guess at the current category
17    guess = classifier.classify(fulltext)
18    entries[x]['pred'] = str(guess[0])
19    entries[x]['cprob'] = round(guess[1], 3)
20    print 'Guess: '+ str(guess[0])
```

Listing 4: Second 50 Classification

In order to create the table, the title, classifier strings, predicted category, actual category, and cprob were sent to "table.txt." The output file was then converted to "table.tex," in order to create table 2.1. Listing 5 is the code to output the data. [4]

```
1  #write results to file, for table creation
2  tab = open('table.txt', 'w', 0)
3  tab.write('Title\tClassifier\tPredicted\tActual\tcprob()\n')
4  for x in range(0, 50):
5    tab.write(entries[x]['title'] + '\t' + entries[x]['classifier'] + '\t' + ' '+ '\t' + entries[x]['actual'
       ]+'\n')
6
7  for x in range(50, 100):
8    tab.write(entries[x]['title'] + '\t' + ' ' + '\t' + entries[x]['pred']+'\t' + entries[x]['actual']+'\t'+
         str(entries[x]['cprob'])+'\n')
9
10 tab.close()
```

Listing 5: Table Data

Table 2.1: Classifier Data

| TITLE | CLASSIFIER | PREDICTED | ACTUAL | cprob() |
|---|---|---|---|---|
| Qubes R3/Odyssey initial source code release | Join us | | announcements | |
| Announcing Qubes OS Release 2! | bug fixes | | announcements | |
| Physical separation vs. Software compartmentalization | topics discussed | | paper | |
| Qubes OS R2 rc2, Debian template, SSLed Wiki, BadUSB, and more... | updates | | announcements | |
| Qubes OS R2 rc1 has been released! | improvements | | announcements | |
| Shattering the myths of Windows security | Windows security | | security | |
| Qubes R2 Beta 3 has been released! | beta | | announcements | |
| Windows 7 seamless GUI integration coming to Qubes OS! | Support Tools | | os | |
| Thoughts on Intel's upcoming Software Guard Extensions (Part 2) | memory bus | | hardware | |
| Thoughts on Intel's upcoming Software Guard Extensions (Part 1) | DRAM | | hardware | |
| Qubes OS R3 Alpha preview: Odyssey HAL in action! | configuration template | | announcements | |
| Introducing Qubes Odyssey Framework | hardcoded | | announcements | |
| Qubes 2 Beta 2 has been released! | drivers | | announcements | |

Table 2.1: Classifier Data

| TITLE | CLASSIFIER | PREDICTED | ACTUAL | cprob() |
|---|---|---|---|---|
| Converting untrusted PDFs into trusted ones: The Qubes Way | PDFs | | paper | |
| Qubes 2 Beta 1 with initial Windows support has been released! | windows support | | os | |
| How is Qubes OS different from... | linux BSD | | os | |
| Introducing Qubes 1.0! | introducing Qubes | | announcements | |
| Qubes 1.0 Release Candidate 1! | release candidate | | announcements | |
| Some comments on "Operation High Roller" | two-factor authentication | | security | |
| Windows support coming to Qubes! | coming to | | announcements | |
| Qubes Beta 3! | installation guide | | announcements | |
| Thoughts on DeepSafe | firewall | | security | |
| Trusted Execution In Untrusted Cloud | confidentiality | | security | |
| Exploring new lands on Intel CPUs (SINIT code execution hijacking) | processors | | hardware | |
| Playing with Qubes Networking for Fun and Profit | infrastructure | | os | |
| Qubes Beta 2 Released! | proud to announce | | announcements | |
| Anti Evil Maid | disk encryption | | security | |
| Interview about Qubes OS | interview about | | announcements | |
| My SSTIC 2011 slides | security | | security | |
| From Slides to Silicon in 3 years! | presentation | | announcements | |
| USB Security Challenges | software attacks | | security | |
| (Un)Trusting the Cloud | encrypted | | security | |
| The App-oriented UI Model and its Security Implications | security by isolation | | security | |
| Following the White Rabbit: Software Attacks Against Intel VT-d | publish | | paper | |
| The Linux Security Circus: On GUI isolation | sniff all | | security | |
| Why the US "password revolution" won't work | multi-factor authentication | | security | |
| Qubes Beta 1 has been released! | new features | | announcements | |
| Partitioning my digital life into security domains | compromise | | security | |
| My documents got lost/stolen [offtopic] | SSL cert | | security | |
| Update on Qubes | Stay tuned | | announcements | |
| Qubes Alpha 3! | milestone | | announcements | |
| ITL is hiring! | looking to hire | | announcements | |
| On Thin Clients Security | desktop security | | security | |
| (Un)Trusting your GUI Subsystem | deprivileged | | security | |
| Qubes, Qubes Pro, and the Future... | the future | | announcements | |
| The MS-DOS Security Model | attacker exploiting | | security | |
| Skeletons Hidden in the Linux Closet: r00ting your Linux Desktop for Fun and Profit | malicious | | security | |
| Qubes Alpha 2 released! | are here | | announcements | |
| Disposable VMs | non-sensitive | | security | |
| On Formally Verified Microkernels (and on attacking them) | verification attempts | | security | |

Table 2.1: Classifier Data

| TITLE | CLASSIFIER | PREDICTED | ACTUAL | cprob() |
|---|---|---|---|---|
| Evolution | | announcements | announcements | 1.0 |
| Remotely Attacking Network Cards (or why we do need VT-d and TXT) | | announcements | security | 1.0 |
| Introducing Qubes OS | | announcements | announcements | 1.0 |
| Priorities | | security | security | 1.0 |
| Another TXT Attack | | announcements | security | 1.0 |
| Evil Maid goes after TrueCrypt! | | announcements | security | 1.0 |
| Intel Security Summit: the slides | | announcements | paper | 1.0 |
| About Apple's Security Foundations Or Lack Of Thereof... | | security | security | 1.0 |
| PDF signing and beyond | | security | security | 1.0 |
| Vegas Toys (Part I): The Ring -3 Tools | | announcements | security | 1.0 |
| Black Hat 2009 Slides | | announcements | paper | 1.0 |
| Interview | | announcements | announcements | 0.983 |
| Virtualization (In)Security Training in Vegas | | announcements | announcements | 1.0 |
| Quest to The Core | | announcements | announcements | 1.0 |
| More Thoughts on CPU backdoors | | security | security | 1.0 |
| Thoughts About Trusted Computing | | announcements | security | 1.0 |
| Trusting Hardware | | security | security | 1.0 |
| The Sky Is Falling? | | security | security | 1.0 |
| Attacking SMM Memory via Intel CPU Cache Poisoning | | announcements | security | 1.0 |
| Independent Attack Discoveries | | announcements | security | 1.0 |
| Attacking Intel TXT: paper and slides | | announcements | paper | 0.997 |
| Nesting VMMs, Reloaded. | | announcements | announcements | 0.997 |
| Closed Source Conspiracy | | security | security | 1.0 |
| Why do I miss Microsoft Bit-Locker? | | security | security | 1.0 |
| Attacking Intel Trusted Execution Technology | | security | security | 1.0 |
| Microsoft executive "rebuts" our research! | | announcements | security | 1.0 |
| Xen 0wning Trilogy: code, demos and q35 attack details posted | | announcements | security | 1.0 |
| The three approaches to computer security | | security | security | 1.0 |
| Teamwork Crediting | | security | security | 1.0 |
| Intel patches the Q35 bug | | announcements | announcements | 1.0 |
| Attacking Xen: DomU vs. Dom0 consideration | | announcements | security | 1.0 |
| Our Xen 0wning Trilogy Highlights | | security | announcements | 1.0 |
| 0wning Xen in Vegas! | | announcements | announcements | 1.0 |
| Rafal Wojtczuk joins Invisible Things Lab | | announcements | announcements | 1.0 |
| 1984? | | announcements | announcements | 0.999 |
| Vegas Training 2008 | | announcements | announcements | 1.0 |
| Research Obfuscated | | security | hardware | 1.0 |

Table 2.1: Classifier Data

| TITLE | CLASSIFIER | PREDICTED | ACTUAL | cprob() |
|---|---|---|---|---|
| The Most Stupid Security News Ever | | security | security | 1.0 |
| The RSA Absurd | | security | security | 1.0 |
| Kick Ass Hypervisor Nesting! | | security | os | 1.0 |
| Razor-Thin Hypervisors | | announcements | announcements | 1.0 |
| Thoughts On Browser Rootkits | | announcements | security | 1.0 |
| Tricky Tricks | | announcements | security | 1.0 |
| Virtualization Detection vs. Blue Pill Detection | | announcements | announcements | 1.0 |
| We're ready for the Ptacek's challenge! | | announcements | announcements | 1.0 |
| Invisible Things Lab, Bitlocker/TPM bypassing and some conference thoughts | | announcements | security | 1.0 |
| Understanding Stealth Malware | | announcements | security | 1.0 |
| The Human Factor | | security | security | 1.0 |
| The Game Is Over! | | security | security | 1.0 |
| Handy Tool To Play with Windows Integrity Levels | | announcements | os | 1.0 |

## 3 PERFORMANCE

Assess the performance of your classifier in each of your categories by computing precision, recall, and F1. Note that the definitions of precisions and recall are slightly different in the context of classification; see:

http://en.wikipedia.org/wiki/Precision_and_recall#Definition_.28classification_context.29
and
http://en.wikipedia.org/wiki/F1_score

### 3.1 SOLUTION

For a classifier, the expected prediction will be positive if labelled correctly and negative if not. The observation will be true or false, as determined by external judgement. A true positive (TP) is an entry labelled correctly. A false negative (FN) would be an entry that was not labelled, but it should have been. This classifier did not have any FN. A false positive (FP) was incorrectly labelled. [2] Precision, or positive prediction value, was calculated as:

$$\text{Precision} = \frac{tp}{tp + fp}$$

The recall, or sensitivity, was calculated as:

$$\text{Recall} = \frac{tp}{tp + fn}$$

The F1 score, or accuracy, was calculated as:

$$F_1 = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

Listing 6 is the code used to perform the calculations. [2], [3]

```
1  tp = 0
2  fn = 0
3  fp = 0
4  for x in range(50, 100):
5    if entries[x]['pred'] == entries[x]['actual']:
6      tp += 1
7    elif not entries[x]['pred'] and entries[x]['actual']:
8      fn += 1
9    else:  #labelled incorrectly
10     fp += 1
11
12 precision = float(tp) / (tp + fp)
13 recall = float(tp) / (tp + fn)
14
15 f1 = 2 * ((precision * recall) / (precision + recall))
```

Listing 6: Calculations

The classifier precision was only .58 and because there were no FN, the recall score was 1. This provided an F1 score of 0.734. I believe this was due to the lack of training for the paper and os categories. There were far fewer entries that belonged here, so they were not classified correctly by the classifier. Perhaps if there had been more trained entries, it would have been higher.

Finally, all entries and data were output to "entries.txt" to be included in this report. [1] Listing 7 in section 4 is the entire classify.py program and listing 8, inthe same section is feedfilter.py. The file docclass.py, as used, is included with this report.

# 4 PYTHON CODE

```
1
2  import urllib
3  import feedparser
4  from bs4 import BeautifulSoup
5  import requests
6  import docclass
7  import feedfilter
8  import re
9
10 entries = []
11 blog = 'http://theinvisiblethings.blogspot.com/'
12 path = '/feeds/posts/default'
13 query_arg = {'alt' : 'atom'}
14 udata = urllib.urlencode(query_arg)
15
16 full_url = blog+path+"?"+udata
17 r = requests.get(full_url)
18 f=feedparser.parse(full_url)
19 entries = f['entries']
20 bsdata = r.text
21 soup = BeautifulSoup(bsdata)
22 next = soup.find('link', rel='next')
23 for x in range(1,4):
24   n= next.get('href')
25   f=feedparser.parse(n)
26   entries.extend(f['entries'])
27   r = requests.get(n)
28   bsdata = r.text
29   soup = BeautifulSoup(bsdata)
30   next = soup.find('link', rel='next')
31
32 _illegal_xml_chars_RE = re.compile(u'[\x00-\x08\x0b\xa9\xae\x0c\x0e-\x1F\uD800-\uDFFF\uFFFE\uFFFF]')
33
34 def remove_tags(text):
35     cleanr =re.compile('<.*?>')
36     cleantext = re.sub(cleanr,'', text)
37     cleantext = re.sub("\n", ' ', cleantext)
38     cleantext = cleantext.strip()
39     cleantext = cleantext.replace(u'\u201c', '')
40     cleantext = cleantext.replace(u'\u201d', '')
41     cleantext = cleantext.replace(u'\u2019', '\'')
42     cleantext = cleantext.replace('-&gt;', '')
43     return _illegal_xml_chars_RE.sub(' ', cleantext)
44
45 for x in range(0,100):
46   entries[x]['content'][0]['value'] = remove_tags(entries[x]['content'][0]['value'])
47   entries[x]['title'] = remove_tags(entries[x]['title'])
48
49 cl=docclass.fisherclassifier(docclass.getwords)
50 cl.setdb('feed.db')
51
52 feedfilter.read(entries, cl)
53
54 #write results to file, for table creation
55 tab = open('table.txt', 'w', 0)
56 tab.write('Title\tClassifier\tPredicted\tActual\tcprob()\n')
57 for x in range(0, 50):
58   tab.write(entries[x]['title'] + '\t' + entries[x]['classifier'] + '\t' + ' '+ '\t' + entries[x]['actual'
        ]+'\n')
59
60 for x in range(50, 100):
61   tab.write(entries[x]['title'] + '\t' + ' ' + '\t' + entries[x]['pred']+'\t' + entries[x]['actual']+'\t'+
        str(entries[x]['cprob'])+'\n')
62
63 tab.close()
64
65 # compute Precision, recall, f1
66 # tp is labelled correctly, fn is not labelled but should have been, fp is incorrectly labelled
67 tp = 0
68 fn = 0
```

```
69  fp = 0
70  for x in range(50, 100):
71    if entries[x]['pred'] == entries[x]['actual']:
72      tp += 1
73    elif not entries[x]['pred'] and entries[x]['actual']:
74      fn += 1
75    else:    #labelled incorrectly
76      fp += 1
77
78  precision = float(tp) / (tp + fp)
79  recall = float(tp) / (tp + fn)
80
81  f1 = 2 * ((precision * recall) / (precision + recall))
82
83  # print the entries to file
84  out = open('entries.txt', 'w', 0)
85  for e in entries:
86    print>>out, e
87
88  out.write('\n---------------------\n')
89  out.write("Precision: " + str(precision)+'\n')
90  out.write("Recall: " + str(recall)+'\n')
91  out.write("F1: " + str(f1)+'\n')
92  out.close()
93
94  print "Precision: " + str(precision)
95  print "Recall: " + str(recall)
96  print "F1: " + str(f1)
```

Listing 7: Complete Classifier

```
1   import feedparser
2   import re
3
4   # Takes a filename of URL of a blog feed and classifies the entries
5   def read(entries, classifier):
6   # classify the first 50 and train the classifer
7     for x in range(0, 50):
8       print
9       print '-----'
10      # Print the contents of the entry
11      print 'Title:      '+entries[x]['title'].encode('utf-8')
12      print
13      print entries[x]['content'][0]['value'].encode('utf-8')
14
15      # Combine all the text to create one item for the classifier
16      fulltext= entries[x]['title'] + ', ' + entries[x]['content'][0]['value']
17
18      # Ask the user to specify the correct category and train on that
19      tr=raw_input('Enter classifier: ')
20      entries[x]['classifier'] = tr
21      cl=raw_input('Enter category: ')
22      entries[x]['actual'] = cl
23      classifier.train(fulltext,cl)
24
25  # classify the final 50 and get the guess and cprob
26      for x in range(50, 100):
27      print
28      print '-----'
29      # Print the contents of the entry
30      print 'Title:      '+entries[x]['title'].encode('utf-8')
31      print
32      print entries[x]['content'][0]['value'].encode('utf-8')
33
34      # Combine all the text to create one item for the classifier
35      fulltext= entries[x]['title'] + ', ' + entries[x]['content'][0]['value']
36
37      # Ask the user to specify the correct category
38      cl=raw_input('Enter category: ')
39      entries[x]['actual'] = cl
40      # Print the best guess at the current category
41      guess = classifier.classify(fulltext)
```

```
42     entries[x]['pred'] = str(guess[0])
43     entries[x]['cprob'] = round(guess[1], 3)
44     print 'Guess: '+ str(guess[0])
```

Listing 8: Modified Feed Filter

# REFERENCES

[1] Alex Martelli. Python: Write a list to a file. `http://stackoverflow.com/questions/899103/python-write-a-list-to-a-file`. Accessed: 2014-12-7.

[2] Precision and recall. `http://en.wikipedia.org/wiki/Precision_and_recall#Definition_.28classification_context.29`. Accessed: 2014-12-10.

[3] F1 score. `http://en.wikipedia.org/wiki/F1_score`. Accessed: 2014-12-10.

[4] Toby Segaran. *Programming Collective Intelligence.* O'Reilly Media, 2007.