# CS 532: Assignment 9

Dinesh Kumar Paladhi

Spring 2016

# Contents

# 1 Problem 1

1. Choose a blog or a newsfeed (or something similar with an Atom
or RSS feed). Every student should do a unique feed, so please
''claim'' the feed on the class email list (first come, first served).
It should be on a topic or topics of which you are qualified to
provide classification training data. Find something with at least
100 entries (or items if RSS).

Create between four and eight different categories for the entries
in the feed:

examples:

work, class, family, news, deals

liberal, conservative, moderate, libertarian

sports, local, financial, national, international, entertainment

metal, electronic, ambient, folk, hip-hop, pop

Download and process the pages of the feed as per the week 12
class slides.

Be sure to upload the raw data (Atom or RSS) to your github account.

## 1.1 Solution

1. I started doing this problem by first searchin for a blog which gives me minimum 100 feeds.

2. It became very tough for me to find a blog like that, Finally I have found this page "http://www.thehindu.com/news/nati
with minimum of 100 rss feeds.

3. There were more than 100 rss feeds, but I picked only 100 from them and saved it in an xml file. Sample
xml file can be seen in fig1.

4. Once I got these feeds I have read all of those titles and description of each feed and then I categorized
them into the following categories.

5. Accident- Anything that is attacked,killed,injured,any loss happened,any illness or any one arrested.
All these situations are categorized into accident category.

6. law- Any anouncement from supreme or high courts,any legal statements are categorized into law
category.

7. politics- Any minister,government and its anouncements and activities are categorized into politics
category.

8. elections- Any votings,campaigns,candidates,polls are categorized into elections category.

9. entertainment- Any anouncement regarding films,any happy moments,movies,songs,celebrations are
categorized into entertainment category.

10. others- All those which do not come in above categories are categorized into others category.

11. transportation- Any information regarding buses,trains,travel details,traffic are categorized into trans-
portation category.

## 1.2 Outputs

**Sample Blog XML file**

```xml
<?xml version="1.0" encoding="UTF-8"?>
<rss version="2.0">
<channel>
<title>The Hindu - National</title>
<link>http://www.thehindu.com/</link>
<description>RSS feed</description>
<language>en-us</language>
<copyright>Copyright 2016 The Hindu</copyright>
<item>
<title><![CDATA[AAP govt. withdraws plea on distribution of powers]]></title>
<author><![CDATA[Legal Correspondent]]></author>
<category><![CDATA[Delhi]]></category>
<link>http://www.thehindu.com/news/cities/Delhi/aap-govt-withdraws-plea-on-distribution-of-powers/article849
<description><![CDATA[
The case is centred around a notification issues by the Centre last year
]]>
</description>
<pubDate><![CDATA[Wed, 20 Apr 2016 02:21:03 +0530]]></pubDate>
</item>
<item>
<title><![CDATA[Uphaar fire tragedy: SC may hear review pleas next month]]></title>
<author><![CDATA[Legal Correspondent]]></author>
<category><![CDATA[Delhi]]></category>
<link>http://www.thehindu.com/news/cities/Delhi/uphaar-fire-tragedy-sc-may-hear-review-pleas-next-month/arti
<description><![CDATA[
Victims' families, CBI have raised questions over the punishment meted out to the accused
]]>
</description>
<pubDate><![CDATA[Wed, 20 Apr 2016 02:19:51 +0530]]></pubDate>
</item>
<item>
```

Figure 1: Sample list of Blog Feeds

# 2 Problem 2

2. Manually classify the first 50 entries, and then classify (using
the fisher classifier) the remaining 50 entries.

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

## 2.1 Solution

1. In order to classify the feeds into the respective categories I need to train the data first.

2. So what I did first is to divide the data into two parts. First 50 feeds are for training the code and the
   next 50 are for getting automatic classifications based on training.

3. This is done with the help of docclass and feedfilter codes which are taken as a reference from Pro-
   gramming collective Intelligence textbook.

4. Later by using the code I figured out to get the cprob() and fisherprob() values for the feeds.

5. I have saved all those outputs and can be seen in the following tables below.

6. Sample output for the first 50 feeds can be seen in the fig2 and the sample output for the last 50 feeds
   can be seen in the fig3/.

## 2.2 Code Listing

```
1
2  import feedparser
3  import re
4  import math
5  import docclass
6  # Takes a filename or URL of a blog feed and classifies the entries
7  def read(feed, classifier):
8
9    splitRegexp = re.compile( r"<[^>]+>" )
10
11   num=0
12   Get feed entries and loop over them
13   f=feedparser.parse(feed)
14   print
15   print '----- Begin manual classification (training) -----'
16   for entry in f['entries'][0:50]:
17     num=num +1
18     # Print the contents of the entry
19     title=entry['title'].encode('utf-8').replace("'","")
20     print 'Title:      '+ title
21
22     description = splitRegexp.sub( "", entry[ "description" ] )
23
24     print description #entry['description'].encode('utf-8')
25
26     # Combine all the text to create one item for the classifier
27     #fulltext='%s\n%s\n%s' % (entry['title'],entry['publisher'],entry['description'])
28     fulltext='%s\n%s' % (entry['title'],entry['description'])
29     # Remove apostrophes
30     fulltext=fulltext.replace("'","")
31     # Print the best guess at the current category
32     predicted=str(classifier.classify(fulltext))
33     print 'Predicted category: ', predicted
34
35     # Ask the user to specify the correct category and train on that
36     actual=raw_input('Actual category: ')
37     feature=None
38     classifier.train(fulltext, actual)
39
40     # Save the manual classifications
41     # num, entry, feature, predicted, actual, cprob=None
42     classifier.manualClassdb(num, title, feature, predicted, actual)
43
44  #def autoClassify(feed, classifier):
45    num=50
46   print '----- Begin automatic classification -----'
47   # Get feed entries and loop over them
48   f=feedparser.parse(feed)
49   for entry in f['entries'][50:100]:
50     num=num+1
51     # Print the contents of the entry
52     title=entry['title'].encode('utf-8').replace("'","")
53     print 'Title:      '+ title
54     description = splitRegexp.sub( "", entry[ "description" ] )
55
56     print description #entry['description'].encode('utf-8')
57
58     # Combine all the text to create one item for the classifier
59     #fulltext='%s\n%s\n%s' % (entry['title'],entry['publisher'],entry['description'])
60     fulltext='%s\n%s' % (entry['title'],entry['description'])
61     fulltext=fulltext.replace("'","")
62     # Print the best guess at the current category
63     predicted=str(classifier.classify(fulltext))
64     print 'Predicted: ', predicted
65
66     # Ask the user to specify the correct category
67     actual=raw_input('Enter actual category: ')
68     feature=raw_input('Enter string classifier: ')
69
70     #classifier.train(entry, cl)
```

```
71        # probability the item should be in this category
72        cp=round( classifier.cprob(feature , predicted ) ,3)
73        print 'cprob: ', str(cp)
74        fischerprob1=round( classifier.fisherprob(feature , predicted ) ,4)
75        print 'fisherprob: ', str(fischerprob1)
76        # Save the trained classifications
77        # num, entry , feature , predicted , actual , cprob(feature , predicted)
78        classifier.autoClassdb(num, title , feature , predicted , actual , cp)
79       # entryfeatures (entry)
80     #return classifier
81
82  def entryfeatures(entry):
83     splitter=re.compile('\\W*')
84     f={}
85
86     # Extract the title words and annotate
87     titlewords=[s.lower() for s in splitter.split(entry['title'])
88            if len(s)>2 and len(s)<20]
89     for w in titlewords: f['Title:'+w]=1
90
91     # Extract the description words
92     descriptionwords=[s.lower() for s in splitter.split(entry['description'])
93            if len(s)>2 and len(s)<20]
94
95     # Count uppercase words
96     uc=0
97     for i in range(len(descriptionwords)):
98       w=descriptionwords[i]
99       f[w]=1
100      if w.isupper(): uc+=1
101
102      # Get word pairs in description as features
103      if i<len(descriptionwords)-1:
104        twowords=' '.join(descriptionwords[i:i+1])
105        f[twowords]=1
106
107    # Removed: Keep creator and publisher whole
108    #f['Publisher:'+entry['publisher']]=1
109
110    # UPPERCASE is a virtual word flagging too much shouting
111    if float(uc)/len(descriptionwords)>0.3: f['UPPERCASE']=1
112    print f
113    return f
114
115 def main():
116    cl=docclass.fisherclassifier(docclass.getwords)
117    cl.setdb('dpaladhi.db')
118    read('my_data.xml',cl)
119 main()
```

Listing 1: Python Code for feedfilter

## 2.3   Code Listing

```
1   #from pysqlite2 import dbapi2 as sqlite
2   from sqlite3 import dbapi2 as sqlite
3   import re
4   import math
5
6   def getwords(doc):
7       splitter=re.compile('\\W*')
8       ## Remove all the HTML tags
9       doc=re.compile(r'<[^>]+>').sub('',doc)
10      # Split the words by non-alpha characters
11      words=[s.lower() for s in splitter.split(doc)
12              if len(s)>2 and len(s)<20]
13
14      # Return the unique set of words only
15      return dict([(w,1) for w in words])
16
17  class classifier:
18      def __init__(self,getfeatures,filename=None):
19          # Counts of feature/category combinations
20          self.fc={}
21          # Counts of documents in each category
22          self.cc={}
23          ## extract features for classification
24          self.getfeatures=getfeatures
25
26      def setdb(self,dbfile):
27          self.con=sqlite.connect(dbfile)
28          self.con.execute('create table if not exists rss(num, entry, feature, predicted, actual,
                  cprob)')
29          self.con.execute('create table if not exists fc(feature,category,count)')
30          self.con.execute('create table if not exists cc(category,count)')
31          # remove old data from previous sessions
32          # self.con.execute('delete from rss')
33          # self.con.execute('delete from fc')
34          # self.con.execute('delete from cc')
35
36      def manualClassdb(self,num, entry, feature, predicted, actual):
37          self.con.execute("insert into rss values ('%s','%s', '%s', '%s','%s', '%s')"
38                          % (num, entry, feature, predicted, actual, None))
39          self.con.commit()
40
41      def autoClassdb(self,num, entry, feature, predicted, actual, cp):
42          self.con.execute("insert into rss values ('%s','%s', '%s', '%s','%s', '%s')"
43                          % (num, entry, feature, predicted, actual, cp))
44          self.con.commit()
45      ## Increase the count of a feature/category pair
46      def incf(self,f,cat):
47          count=self.fcount(f,cat)
48          if count==0:
49              self.con.execute("insert into fc values ('%s','%s',1)"
50                              % (f,cat.lower()))
51          else:
52              self.con.execute(
53                  "update fc set count=%d where feature='%s' and category='%s'"
54                  % (count+1,f,cat.lower()))
55
56      ## The number of times a feature has appeared in a category
57      def fcount(self,f,cat):
58          res=self.con.execute(
59              'select count from fc where feature="%s" and category="%s"'
60              %(f,cat)).fetchone()
61          if res==None: return 0
62          else: return float(res[0])
63
64      ## Increase the count of a category
65      def incc(self,cat):
66          count=self.catcount(cat)
67          if count==0:
68              self.con.execute("insert into cc values ('%s',1)" % (cat.lower()))
69          else:
```

```python
        self.con.execute("update cc set count=%d where category='%s'"
                         % (count+1,cat))

    ## The number of items in a category
    def catcount(self,cat):
        res=self.con.execute('select count from cc where category="%s"'
                             %(cat)).fetchone()
        if res==None: return 0
        else: return float(res[0])

    ## The list of all categories
    def categories(self):
        cur=self.con.execute('select category from cc');
        return [d[0] for d in cur]

    ## The total number of items
    def totalcount(self):
        res=self.con.execute('select sum(count) from cc').fetchone();
        if res==None: return 0
        return res[0]


    ## The train method takes an item(document) and a classification.
    ## It uses the getfeatures function to the break the item into its
    ## separate features. It then calls incf to increase the counts for
    ## this classification for every feature. Finally, it increases
    ## the total count for this classification.
    def train(self,item,cat):
        features=self.getfeatures(item)
        # Increment the count for every feature with this category
        for f in features:
            self.incf(f,cat)

        # Increment the count for this category
        self.incc(cat)
        self.con.commit()

    ## Probability is a number between 0 and 1, indicating
    ## the likelihood of an event. You calculate the probability of
    ## a word in a particular category by dividing the number of
    ## times the word appears in a document in that category
    ## by the total number of documents in the category.
    def fprob(self,f,cat):
        if self.catcount(cat)==0: return 0

        # The total number of times this feature appeared in this
        # category divided by the total number of items in this category
        return self.fcount(f,cat)/self.catcount(cat)

    def weightedprob(self,f,cat,prf,weight=1.0,ap=0.5):
        # Calculate current probability
        basicprob=prf(f,cat)

        # Count the number of times this feature has appeared in
        # all categories
        totals=sum([self.fcount(f,c) for c in self.categories()])

        # Calculate the weighted average
        bp=((weight*ap)+(totals*basicprob))/(weight+totals)
        return bp




class naivebayes(classifier):

    def __init__(self,getfeatures):
        classifier.__init__(self,getfeatures)
        self.thresholds={}

    def docprob(self,item,cat):
        features=self.getfeatures(item)
```

```python
142
143         # Multiply the probabilities of all the features together
144         p=1
145         for f in features: p*=self.weightedprob(f,cat,self.fprob)
146         return p
147
148     def prob(self,item,cat):
149         catprob=self.catcount(cat)/self.totalcount()
150         docprob=self.docprob(item,cat)
151         return docprob*catprob
152
153     def setthreshold(self,cat,t):
154         self.thresholds[cat]=t
155
156     def getthreshold(self,cat):
157         if cat not in self.thresholds: return 1.0
158         return self.thresholds[cat]
159
160     def classify(self,item,default=None):
161         probs={}
162         # Find the category with the highest probability
163         max=0.0
164         for cat in self.categories():
165             probs[cat]=self.prob(item,cat)
166             if probs[cat]>max:
167                 max=probs[cat]
168                 best=cat
169
170         # Make sure the probability exceeds threshold*next best
171         for cat in probs:
172             if cat==best: continue
173             if probs[cat]*self.getthreshold(best)>probs[best]: return default
174         return best
175
176 ## This function will return the probability that an item with the
177 ## specified feature belongs in the specified category, assuming there
178 ## will be an equal number of items in each category.
179 class fisherclassifier(classifier):
180     def cprob(self,f,cat):
181         # The frequency of this feature in this category
182         clf=self.fprob(f,cat)
183         if clf==0: return 0
184
185         # The frequency of this feature in all the categories
186         freqsum=sum([self.fprob(f,c) for c in self.categories()])
187
188         # The probability is the frequency in this category divided by
189         # the overall frequency
190         p=clf/(freqsum)
191
192         return p
193
194
195     def fisherprob(self,item,cat):
196         # Multiply all the probabilities together
197         p=1
198         features=self.getfeatures(item)
199         for f in features:
200             p*=(self.weightedprob(f,cat,self.cprob))
201
202         # Take the natural log and multiply by -2
203         fscore=-2*math.log(p)
204         fprobvalue=self.invchi2(fscore,len(features)*2)
205         #print fprobvalue
206
207         # Use the inverse chi2 function to get a probability
208         return fprobvalue
209
210     ## Inverse chi-squared function
211     def invchi2(self,chi, df):
212         m = chi / 2.0
213         sum = term = math.exp(-m)
```

```
214        for i in range(1, df//2):
215            term *= m / i
216            sum += term
217        return min(sum, 1.0)
218
219    def __init__(self, getfeatures):
220        classifier.__init__(self, getfeatures)
221        self.minimums={}
222
223    def setminimum(self, cat, min):
224        self.minimums[cat]=min
225
226    def getminimum(self, cat):
227        if cat not in self.minimums: return 0
228        return self.minimums[cat]
229
230    def classify(self, item, default=None):
231        # Loop through looking for the best result
232        best=default
233        max=0.0
234        for c in self.categories():
235            p=self.fisherprob(item,c)
236            # Make sure it exceeds its minimum
237            if p>self.getminimum(c) and p>max:
238                best=c
239                max=p
240        return best
```

Listing 2: Python Code for docclass

## 2.4 Outigitik

**Output 1**



```
Title:      Youth turn bird saviours in Bidar
Two young members of the Bidar Photographic Society (BPS) are inspiring others by putting up pots of gra
Predicted category:  others
Actual category: others
Title:      Nagaland Chief Secretary is first Ambassador for Girl Child project
Nagaland Chief Secretary Pankaj Kumar has become the first Ambassador for Girl Child (AFGC) under the Ce
Kumar also released the ...
Predicted category:  politics
Actual category: politics
Title:      Lockdown at Peenya Industrial Area
Trade bodies suspect the losses will run into several crores
Predicted category:  politics
Actual category: others
Title:      Timely action by police, fire services prevents fire mishap
An 18-tonne LPG bullet tanker overturned near Kalladka, 30 km from Mangaluru, on the Mangaluru- Bengalu
Predicted category:  accident
Actual category: accident
Title:      AIUTUC condemns lathicharge on garment workers
The All-India United Trade Union Centre (AIUTUC) has strongly condemned Monday's police lathicharge in B
Predicted category:  politics
Actual category: accident
Title:      Ghulam Ali's performance in Bhavnagar cancelled
Noted Pakistani ghazal singer Ghulam Ali, who was scheduled to perform on Tuesday at a cultural event in
Predicted category:  accident
Actual category: entertainment
Title:      Bengaluru blockade: KSRTC suspends Mysuru-Bengaluru bus service
KSRTC has suspended its services from Mysuru to Bengaluru in view of the road blockade between Ramanaga
Though buses that left Mysuru early in the morning around 5.30 ...
Predicted category:  transportation
Actual category: transportation
Title:      Three killed due to asphyxiation
Three youths died due to asphyxiation while cleaning an old well in Nadanga village in Sirguppa taluk o
Predicted category:  accident
Actual category: accident
Title:      Traffic on Mysuru-Bengaluru highway hit
Police divert vehicles through smaller roads
Predicted category:  transportation
Actual category: transportation
Title:      Most hitches in Rafale deal addressed: Govt.
Most of the hitches in the negotiations with France for the direct purchase of 36 Rafale fighter jets ha
Predicted category:  politics
Actual category: politics
```

Figure 2: Sample outputs for first 50 feeds

| Title | Predicted | Actual |
|---|---|---|
| AAP govt. withdraws plea on distribution of powers | none | politics |
| Uphaar fire tragedy: SC may hear review pleas next month | politics | law |
| Where the mind is not without fear | politics | elections |
| Buddhadeb on poll circuit | politics | elections |
| If we win, there will be common programme | elections | elections |
| Nitish Kumar will be next PM: Lalu | elections | elections |
| Union govt. allots Rs. 800 crore to clean up polluted lakes in garden | elections | politics |
| Will get Kohinoor back, says Centre | elections | politics |
| All members of mob equally guilty: HC | elections | law |
| Govt. yet to pay farmers in Punjab for procured wheat | politics | politics |
| SC asks NAAC to hear grievances of deemed varsities | elections | law |
| Developed countries must tax coal for climate fund | politics | politics |
| Kirpals body arrives in India | politics | accident |
| Centre introducing chaos: HC | politics | politics |
| Trading bloc to India: Cut tariffs or exit FTA talks | politics | others |
| PM invokes Vajpayee, moots development of Kashmir | politics | politics |
| Handwara rejoices as Army bunkers are dismantled | politics | accident |
| Two youths create oases for birds as Bidar sizzles | politics | others |
| Centre rejects T.N. proposal to free Rajiv Gandhi killers | politics | politics |
| Rajans choice of words could have been better | law | politics |
| As an alumnus, I feel hurt over JNU controversy: Nirmala | politics | politics |
| Woman Maoist killed in Gadhchiroli encounter | accident | accident |
| Insult to God to have unauthorised places of worship: SC | law | law |
| I have not been formally approached for Atulya Bharat: Amitabh Bachchan | politics | entertainment |
| Odisha to provide free drinking water to urban poor | politics | politics |

Table 1: Manually classified first 25 entries

| Title | Predicted | Actual |
|---|---|---|
| India to insist written commitment from Pak on NIA team visit | entertainment | accident |
| TN diocese sued for reinstating convicted priest | politics | accident |
| Massive effort to be launched for water conservation: Modi | politics | politics |
| Mahaveer Jayanti celebrated | politics | entertainment |
| Forty five more fire stations to be set up in Odisha: Patnaik | politics | others |
| LDF promises to free Vigilance, create 25 lakh jobs | politics | politics |
| Movement of buses affected between Tumakuru, Bengaluru | politics | transportation |
| Explaining Ola and Ubers surge pricing | elections | transportation |
| Trains between Mysuru, Bengaluru packed to capacity | transportation | transportation |
| Tirupati temple deposits 1,311 kg gold in bank | politics | others |
| Mob ransacks Hebbagodi police station | entertainment | accident |
| EC silent on complaints against AIADMK: Pon Radhakrishnan | politics | politics |
| Kalaburagi, Bidar record highest temperature | others | others |
| Women-only bus service launched in Kashmir | entertainment | transportation |
| Goa Government to bring monkey-hunting tribe to mainstream | politics | politics |
| Youth turn bird saviours in Bidar | others | others |
| Nagaland Chief Secretary is first Ambassador for Girl Child project | politics | politics |
| Lockdown at Peenya Industrial Area | politics | others |
| Timely action by police, fire services prevents fire mishap | accident | accident |
| AIUTUC condemns lathicharge on garment workers | politics | accident |
| Ghulam Alis performance in Bhavnagar cancelled | accident | entertainment |
| Bengaluru blockade: KSRTC suspends Mysuru-Bengaluru bus service | transportation | transportation |
| Three killed due to asphyxiation | accident | accident |
| Traffic on Mysuru-Bengaluru highway hit | transportation | transportation |
| Most hitches in Rafale deal addressed: Govt. | politics | politics |

Table 2: Manually classified next 25 entries from 25 to 50 entries

| Title | Predicted | Actual | Cprob | fisherprob |
|---|---|---|---|---|
| MBBS applications to be issued from May 9 | others | others | 0.0 | 0.5 |
| Sripad Naik admitted to hospital, discharged | politics | politics | 0.0 | 0.4874 |
| Delhi govt impounds 18 taxis for overcharging | transportation | transportation | 1.0 | 0.8333 |
| Garment workers stir continues in Bengaluru, traffic hit for second day | transportation | transportation | 0.0 | 0.8333 |
| Global economic situation grim, worrisome: Jaitley | transportation | others | 0.0 | 0.5 |
| Pak troops violate ceasefire | politics | politics | 0.0 | 0.7334 |
| Local youth held in Handwara molestation case | accident | accident | 0.0 | 0.5 |
| Adopt Periyars Self-Respect principles, Bhagwan tells deprived classes | politics | accident | 0.0 | 0.5 |
| Unstable academic calendar has made students life messy | politics | politics | 0.0 | 0.9 |
| Cops told to be on alert | entertainment | others | 0.0 | 0.5 |
| Children need protection from heat and dehydration | others | accident | 0.0 | 0.5 |
| Bengaluru Today | transportation | others | 0.0 | 0.5 |
| Waive farm loans even if it means borrowing advance | politics | politics | 0.0 | 0.9 |
| Kabini backwaters, a paradise for animals during drought | others | others | 0.0 | 0.5 |
| New pay still on paper for gazetted officers | others | others | 0.0 | 0.5 |
| Four killed in accident | accident | accident | 0.0 | 0.5 |
| Case registered against polytechnic staff for beating student | politics | accident | 0.0 | 0.5 |
| Training programme on tilapia fish farming | politics | others | 0.0 | 0.5966 |
| PEW busts spurious liquor-manufacturing unit | accident | accident | 0.0 | 0.5 |
| Tremors rock Andaman islands | entertainment | accident | 0.0 | 0.5 |
| JD(S) wins Hassan local body by-election | entertainment | elections | 0.0 | 0.25 |
| Where the grass is green even in blazing summer | transportation | others | 0.0 | 0.5 |
| Tamil Nadu Assembly elections Poll diary | accident | elections | 0.0 | 0.5 |
| Vasan to start his campaign from Papanasam | transportation | elections | 0.0 | 0.25 |
| An occasion to build new relationships | entertainment | entertainment | 0.0 | 0.5 |

Table 3: Automatically classified 25 entries from 50 to 75 entries

| Title | Predicted | Actual | Cprob | fisherprob |
|---|---|---|---|---|
| Bears nocturnal adventure triggers beehive of activity | accident | accident | 0.0 | 0.5 |
| CCTVs to be installed in parts of Karimnagar | entertainment | politics | 0.0 | 0.25 |
| Balineni scoffs at rumours | elections | politics | 0.0 | 0.752 |
| Monkey dies after being attacked by dog | transportation | transportation | 0.0 | 0.8333 |
| Troubled by tradition | politics | politics | 0.0 | 0.5966 |
| Red rebels kill Odisha villager | accident | accident | 0.0 | 0.5 |
| Punjab CM approves Rs 750 crore for roads | politics | politics | 0.0 | 0.3289 |
| 14 cases reported | entertainment | elections | 0.0 | 0.25 |
| Dolphin washed ashore in Kilakkarai | politics | accident | 0.0 | 0.5966 |
| Fill vacancies in High Court | politics | law | 0.0 | 0.2904 |
| engagements | law | entertainment | 0.0 | 0.5 |
| Remembered only during elections, they harbour no high hopes | others | others | 0.0 | 0.6407 |
| Three-time MLA to contest in Cumbum | accident | elections | 0.0 | 0.5 |
| TMC fields candidates for Vilathikulam, Srivaikuntam | transportation | elections | 0.0 | 0.5 |
| High Court reserves order on Virbhadras children plea | politics | law | 0.0 | 0.2904 |
| Are there no rights violations in Union Territories, SC asks Centre | politics | politics | 0.0 | 0.728 |
| Animation film to promote brand Amaravati | others | entertainment | 0.0 | 0.5 |
| Seeking divine help to garner votes | politics | politics | 0.0 | 0.75 |
| Kochis public transport to take new route | transportation | transportation | 1.0 | 0.8333 |
| Ragi gruel centre launched for traffic police | politics | transportation | 0.0 | 0.1667 |
| Cinema | entertainment | entertainment | 0.0 | 0.75 |
| Engagements | accident | entertainment | 0.0 | 0.5 |
| Martyrs chronicle: bringing together bits and pieces | accident | accident | 0.0 | 0.5 |
| Vasan to start his campaign from Papanasam | accident | accident | 0.0 | 0.5 |
| Nine arrested on murder charge | accident | accident | 1.0 | 0.8333 |

Table 4: Automatically classified next 25 entries from 75 to 100 entries

**Output 2**

```
----- Begin automatic classification -----
Title:     MBBS applications to be issued from May 9
A total of 2655 seats are available in 20 government medical colleges.
Predicted:  others
Enter actual category: others
Enter string classifier: colleges
cprob:  0.0
fisherprob:  0.5
Title:     Sripad Naik admitted to hospital, discharged
Union Minister of State for AYUSH(Independent charge) Sripad Naik who was admitt
ed on Tuesday morning to Sub-District hospital in Ponda in south Goa after he co
mplained of neck-pain and high blood p...
Predicted:  politics
Enter actual category: politics
Enter string classifier: Minister
cprob:  0.0
fisherprob:  0.4874
Title:     Delhi govt impounds 18 taxis for overcharging
A day after it warned app-based cab companies Ola and Uber against charging cust
omers more than State-prescribed fares, 18 vehicles were impounded by the Transp
ort Department here on Tuesday.
A sou...
Predicted:  transportation
Enter actual category: transportation
Enter string classifier: transport
cprob:  1.0
fisherprob:  0.8333
Title:     Garment workers stir continues in Bengaluru, traffic hit for second d
ay
Traffic from Bommanahalli and Hosur Road Junction diverted to adjacent areas; pr
otests reported at Yeshwanthpur, Gorguntepalya.
Predicted:  transportation
Enter actual category: transportation
Enter string classifier: Traffic
cprob:  0.0
fisherprob:  0.8333
Title:     Modis convocation address at Shri Mata Vaishno Devi University
The Prime Minister is visiting visit Katra, in Jammu and Kashmir, on Tuesday. He
 will inaugurate the Shri Mata Vaishno Devi Narayana Superspeciality Hospital. H
e will deliver the Convocation Address ...
Predicted:  politics
Enter actual category: politics
Enter string classifier: Prime Minister
cprob:  0.0
fisherprob:  0.7334
```

Figure 3: Sample output for last 50 feeds

# 3 Problem 3

3. Assess the performance of your classifier in each of your
categories by computing precision, recall, and F-measure.

## 3.1 Solution

1. In this I need to calculate the precision, recall and F-measure for each of the category.

2. For doing this I have used the following formula.

$$Precesion = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F - Measure = \frac{2TP}{2TP + FP + FN}$$

3. TN is the true negative that is there is the respective category is not present in predicted and actual fields.

4. TP is the true positive that is there is the respective category is present in both predicted and actual fields.

5. FN is the False negative that is there is the respective category is present in actual and not present in predicted field.

6. FP is the False positive that is there is the respective category is not actual and is present in predicted field.

| Category | TN | TP | FN | FP |
|----------|----|----|----|----|
| accident | 35 | 7 | 5 | 3 |
| law | 47 | 0 | 2 | 1 |
| politics | 33 | 9 | 1 | 7 |
| elections | 43 | 0 | 6 | 1 |
| entertainment | 4 | 2 | 3 | 5 |
| others | 37 | 4 | 7 | 2 |
| transportation | 40 | 3 | 1 | 6 |

Table 5: TN, TP, FN, FP values for each category

| Category | Precision | Recall | F-Measure |
|----------|-----------|--------|-----------|
| accident | 0.7 | 0.5833 | 0.6363 |
| law | 0 | 0 | 0 |
| politics | 0.5625 | 0.9 | 0.6923 |
| elections | 0 | 0 | 0 |
| entertainment | 0.2857 | 0.4 | 0.3333 |
| others | 0.6667 | 0.3636 | 0.4705 |
| transportation | 0.3333 | 0.75 | 0.4615 |

Table 6: Precision, Recall and F-Measure values for each category

# Bibliography

[1] The Hindu, National, http://www.thehindu.com/news/national/, 2016

[2] Wikipedia,Precision and recall,https://en.wikipedia.org/wiki/Precision_and_recall, 2016

[3] Wikipedia, F1 score, https://en.wikipedia.org/wiki/F1_score, 2016