

INTRODUCTION TO WEB SCIENCES: Assignment 9

Manoj Chandra Kompalli

21 April 2016

Contents

1	Question 1:	2
1.1	Approach	2
1.2	Input blog data	3
2	Question 2:	4
2.1	Approach	4
2.2	Tables	4
2.2.1	Manual Classification	4
2.2.2	Last 50 items using fisher classifier	6
2.3	Code Listing	9
2.3.1	feedfilter.py	9
2.3.2	prog.py	12
2.3.3	docclass.py	12
3	Question 3:	16
3.1	Approach	16
3.2	Tables	16

1 Question 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).

1.1 Approach

A blog with more than 100 feeds was really a huge task. Most blogs had around 30 feeds. I had actually tried using a blog about astronomy. I could not categorize the items easily in that blog because most of the items were related.

- I have used a newspaper blogs sports page.
- The Newspaper is Indian Express. As it is an Indian Newspaper, all the sports articles are about the sports followed in India.
- Cricket is big in India and most articles are based on it.
- Football is a world sport so I found most articles about it.
- As the Rio Olympics are approaching, some data was about it.
- Tennis is also a big sport and currently many events are taking place all over the world.
- Hockey is India's national sport and so some articles could be of that.
- For sports like WWE, Motorsports I have the category Others.

1.2 Input blog data

This XML file does not appear to have any style information associated with it. The document tree

```
<rss xmlns:content="http://purl.org/rss/1.0/modules/content/" xmlns:wfw="http://wellformedweb.org/
<channel>
<title>The Indian ExpressSports - The Indian Express</title>
<atom:link href="http://indianexpress.com/section/sports/feed/" rel="self" type="application/rss+xml">
<link>http://indianexpress.com</link>
<description>
Latest News, Breaking News Live, Current Headlines, India News Online
</description>
<lastBuildDate>Thu, 21 Apr 2016 21:09:15 +0000</lastBuildDate>
<language>en</language>
<sy:updatePeriod>hourly</sy:updatePeriod>
<sy:updateFrequency>1</sy:updateFrequency>
<generator>http://wordpress.com/</generator>
<cloud domain="indianexpress.com" port="80" path="/?rsscloud=notify" registerProcedure="" protocol="http">
<image>
<url>
http://0.gravatar.com/blavatar/efe0300e7f891c5c802ed340f6b20b67?s=96&d=http%3A%2F%2Fs2.wp.com%2Fimages%2Favatars%2Favatar_96.jpg
</url>
<title>Sports - The Indian Express</title>
<link>http://indianexpress.com</link>
</image>
<atom:link rel="search" type="application/opensearchdescription+xml" href="http://indianexpress.com/search/" />
<atom:link rel="hub" href="http://indianexpress.com/?pushpress=hub"/>
<item>
<title>
Pune boy to appear in world in-line hockey Championship at New Zealand
</title>
<link>
http://indianexpress.com/article/sports/hockey/iihf-pune-boy-to-appear-in-world-in-line-hockey-championship-at-new-zealand/
</link>
<comments>
http://indianexpress.com/article/sports/hockey/iihf-pune-boy-to-appear-in-world-in-line-hockey-championship-at-new-zealand/
</comments>
<pubDate>Fri, 22 Apr 2016 01:40:06 +0000</pubDate>
<lastBuildDate>Fri, 22 Apr 2016 01:40:06 +0000</lastBuildDate>
<dc:creator>
<![CDATA[ Samreen Sayyed ]]>
</dc:creator>
```

Figure 1: Blog data in xml file containing 100 items

2 Question 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 Approach

I have started off by training the first 50 words. I have used methods like classify, fisherprob, cprob from docclass.py and feedfilter.py . Initially most predictions were incorrect but as the training progressed there were good number of hits mostly in categories like cricket and football.

- I have extracted fisherprob and cprob using the functions from docclass.
- Listings show Fisher method used which predicts category based on the entry.
- Based on the first 50 tarained entries, next 50 categories were predicted.
- I am showing title,actual,predicted for the first 50 feeds.
- I am displaying Title,actual,predicted,,cprob,fprob for the last 50 blog items

2.2 Tables

2.2.1 Manual Classification

No.	Title	Predicted	Actual
1	Pune boy to appear in world in-line hockey Championship at New Zealand	None	hockey
2	Leicester City and Spurs dominate Players of the Year selection	hockey	football
3	Chris Gayle readies to unwind after birth of daughter Blush	hockey	cricket
4	Virat Kohli, Ishant Sharma, Ajinkya Rahane, Irfan Pathan indulge in charitable work	hockey	cricket
5	Leicester Citys Jamie Vardy accepts FA charge, awaits decision on ban	football	football
6	Rafael Nadal beats Albert Montanes to reach Barcelona Open quarterfinals	hockey	tennis
7	Louis Oosthuizen joins Vijay Singh, Adam Scott to give Rio Olympics golf a miss	hockey	olympics
8	IPL 2016: Sunrisers Hyderabad tame the Lions in their own backyard	football	cricket
9	It took guts to pick Marcus Rashford, says Manchester United boss Louis Van Gaal	football	football
3	Chris Gayle readies to unwind after birth of daughter Blush	hockey	cricket

4	Virat Kohli, Ishant Sharma, Ajinkya Rahane, Irfan Pathan indulge in charitable work	hockey	cricket
5	Leicester Citys Jamie Vardy accepts FA charge, awaits decision on ban	football	football
6	Rafael Nadal beats Albert Montanes to reach Barcelona Open quarterfinals	hockey	tennis
7	Louis Oosthuizen joins Vijay Singh, Adam Scott to give Rio Olympics golf a miss	hockey	olympics
8	IPL 2016: Sunrisers Hyderabad tame the Lions in their own backyard	football	cricket
9	It took guts to pick Marcus Rashford, says Manchester United boss Louis Van Gaal	football	football
17	IPL 2016: RCB replace injured Samuel Badree with uncapped South African Tabraiz Shamsi	cricket	cricket
18	Gagan Narang, Chain Singh to vie for final spot in Mens 50m rifle prone	cricket	others
19	Cristiano Ronaldo downplays injury, says all good	cricket	football
20	Jurgen Klopp formula starting to work magic for Liverpool, says James Milner	cricket	football
21	Uncle Virat Kohli congratulates Chris Gayle on birth of daughter Blush	cricket	cricket
22	Rio Games count down starts with Olympia torch lighting	cricket	olympics
23	Ive slept with more than 500 women while on tours for West Indies, former seamer Tino Best reveals in upcoming book	hockey	cricket
24	We, humans, failed you: Virat Kohli tweets on Shaktimans death	cricket	cricket
25	I just want to win, says Englands Ben Stokes	cricket	cricket
26	Leicester City fairytale sprinkles magic dust on rival fans	football	football
27	Jamie Vardy was very unlucky: Roy Hodgson	football	football
28	Chyna, former WWE womens champ, passes away	cricket	others
29	Liverpool thrash 10-man Everton 4-0 in Merseyside derby	football	football
30	Neymar to play for Brazil at Rio Games 2016, not Copa America	cricket	football
31	IPL 2016: KXIP opt for Dharamsala as second home venue, leave BCCI, HPCA confused	cricket	cricket
32	Will look to continue creating history: Dipa Karmakar	cricket	others
33	IPL 2016: Krunal Pandya knows his bowling well, says Rohit Sharma	cricket	cricket
34	Chris Gayle becomes father to girl Blush	cricket	cricket
35	Cristiano Ronaldo needs more rest after injury scare, says Zinedine Zidane	football	football
36	Matteo Darmian shines as Manchester United beat Crystal Palace 2-0	football	football

37	ICC to discuss structure, scheduling of bilateral series	cricket	cricket
38	Madras High Court dismisses PIL against BCCI	cricket	cricket
39	Inzamam-ul-Haq took up Pakistan chief selector job on condition of time to prepare squad for World Cup 2019	cricket	cricket
40	Barcelona snap losing streak to thump Deportivo 8-0; Atletico, Real Madrid win to keep league race alive	cricket	football
41	IPL 2016: Water woes brew in Jaipur	cricket	cricket
42	Chip off old block: Rahul Dravids son Samit scores 125 for school team	cricket	cricket
43	FIFA U-17 World Cup tournament director buoyed by Narendra Modis encouragement	cricket	football
44	Jack Wilshere will be available after West Brom, says Arsene Wenger	cricket	football
45	Bayern Munich keeper Manuel Neuer signs contract extension to 2021	cricket	football
46	Jose Mourinho, Claudio Ranieri to manage at Old Trafford for charity	football	football
47	Usain Bolts main goal: to defend the titles, do a three-peat at Rio Olympics 2016	cricket	olympics
48	I want to go and get more goals for Tottenham Hotspur, says Harry Kane	football	football
49	Bombay High Court allows BCCI to hold May 1 IPL 2016 MI vs RPS match in Pune	cricket	cricket
50	Athletics report: Tourist hurdler, tired athlete and a bumbling walker	cricket	others

Table 1: Entries classified manually

2.2.2 Last 50 items using fisher classifier

No.	Title	Predicted	Actual	CProb	FisherProb
51	Have to improve my landings, says Dipa Karmakar	cricket	Olympics	0	0.083
52	Sports federations concerned about venues for Rio Games 2016	cricket	Olympics	0	0.25
53	KKR put KXIP back in their place bottom	cricket	cricket	0	0.5
54	We have to win, we need to win, says Louis Van Gaal	football	football	0	0.833
55	PSG silent over Jose Mourinho reports	football	football	0	0.5
56	PM Modi lauds Dipa Karmakar for her determination	football	football	0	0.5
57	Pablo Zabaleta casts doubt over Manchester City future	football	football	0	0.5

58	From what Ive heard Antonio Contes disciplined, says Cesc Fabregas	football	football	0	0.5
59	Former German FA boss cleared over Qatar cancer comment	football	football	0	1
60	Formula 1 should stop tinkering with rules: Mercedes boss Toto Wolff	cricket	others	0	0.5
61	Vincent Kompany is working 100 percent with normality, says Manuel Pellegrini	football	football	0	0.75
62	The performance Tottenham Hotspur showed was perfect, says Mauricio Pochettino	football	football	0	0.75
63	SRH vs MI, IPL 2016: David Warner played a great innings, says Tim Southee	cricket	cricket	0	0.955
64	Next objective is a medal in Rio Olympics 2016, says Dipa Karmakar	others	olympics	0	0.083
65	IPL 2016, SRH vs MI: Barinder Sran fined for inappropriate conduct	cricket	cricket	0	1
66	Ex-CBI chief moves High Court over BCCI-ICC revenue model	cricket	cricket	0	0.75
67	Suresh Raina, wife Priyanka expecting first child	cricket	cricket	0	0.75
68	PCB disappointed after West Indies refuse to tour Pakistan in 2016	cricket	cricket	0	0.75
69	South Africa players against day-night Tests, feel disadvantaged	cricket	cricket	0	0.5
70	Spurs hit Stoke City for a four, cut Leicester Citys lead	football	football	0	0.875
71	Novak Djokovic, Serena Williams wins Laureus World Sportsman and Sportswoman of the Year awards	football	tennis	0	0.5
72	Formula One allows Pirelli more track days to test tyres for 2017	cricket	others	0	0.25
73	After win in Monte Carlo, Rafael Nadal turns sights to Barcelona	cricket	tennis	0	0.5
74	At Rio Olympics, India plans to promote Make in India initiative	olympics	olympics	0	0.87

75	One-state, one-vote in BCCI will lead to internal politics, BCA tells Supreme Court	cricket	cricket	0	0.9
76	Dipa Karmakar vaults into history books after qualifying for Rio Olympic Games 2016	others	Olympics	0	0.25
77	IPL 2016, SRH vs MI: In David vs Goliath, Warner the difference	cricket	cricket	0	0.75
78	Charged with improper conduct, Jamie Vardy could face extended ban	football	footbal	0	0.5
79	Newcastle United hoping Manchester City will be distracted by Champions League semi-final	football	football	0	0.5
80	Barcelonas slump baffles Carles Puyol, Luis Figo, Raul	football	football	0	0.5
81	Azlan Shah Cup: A few take-aways for India mens hockey team	cricket	hockey	0	0.25
82	IPL 2016, KXIP vs RPS: All is well that Glenn Maxwell ends	cricket	cricket	0	0.5
83	IPL 2016: Who said what about DDs win over RCB	cricket	cricket	0	0.833
84	Bengaluru FC retain I-League title with one game to play	football	football	0	1
85	Rafael Nadal is King of Clay once again with Monte Carlo title	football	tennis	0	0.25
86	Leaving Afghanistans coach job costs Inzamam-ul-Haq rupees 4 lakhs	cricket	cricket	0	0.75
87	Leicester City, minus Jamie Vardy, salvage 2-2 draw with West Ham United	football	football	0	0.666
88	Jurgen Klopp makes 10 changes as Liverpool beat Bournemouth	football	football	0	0.833
89	IPL 2016, RCB vs DD: DD beat RCB by 7 wickets	cricket	cricket	0	0.5
90	Rajkot police orders inquiry into celebratory firing at Ravindra Jadejas wedding	cricket	cricket	0	0.5
91	Roy Hodgson plays down Andy Carrolls Euro 2016 chances	football	football	0	0.5
92	Sebastian Vettel hits out at torpedo Daniil Kvyat	others	others	0	0.5
93	Engine problems force Lewis Hamilton to back of China grid	football	other	0	0.5

94	Positive start by Indian shooters in Rio World Cup	others	others	0	0.5
95	Harbhajan Singh can take his villa whenever he wants, says Amrapali CMD	football	cricket	0	0.125
96	IPL 2016: KKR beat SRH by eight wicket, Gambhir smashes unbeaten 90	cricket	cricket	0	0.5
97	Mitchell Starc ties the knot with girlfriend Alyssa Healy	cricket	cricket	0	0.5
98	IPL 2016 preview: Struggling MI face RCB challenge at Wankhede Stadium	cricket	cricket	1	0.955
99	IPL 2016, KXIP vs KKR: KXIP seek home comfort against KKR	cricket	cricket	1	0.833
100	IPL 2016, RCB vs DD: Quintal de Kock gives Delhi Daredevils heavyweight scalp	cricket	cricket	1	0.5

Table 2: Entries classified using Fisher Classifier

2.3 Code Listing

2.3.1 feedfilter.py

```

1 import docclass
2 import feedparser
3 import re
4 import math
5
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         summary = splitRegexp.sub( "", entry[ "summary" ] )
23
24         print summary #entry['summary'].encode('utf-8')
25
26
27

```

```

28
29
30
31 # Combine all the text to create one item for the classifier
32 #fulltext='%s\n%s\n%s' % (entry['title'],entry['publisher'],entry['summary'])
33
34 fulltext='%s\n%s' % (entry['title'],entry['summary'])
35
36 # Remove apostrophes
37
38 fulltext=fulltext.replace("'",'')
39
40 # Print the best guess at the current category
41
42 predicted=str(classifier.classify(fulltext))
43 print 'Predicted category: ', predicted
44
45 # Ask the user to specify the correct category and train on that
46
47 actual=raw_input('Actual category: ')
48 feature=None
49 classifier.train(fulltext, actual)
50
51 # Save the manual classifications
52 # num, entry, feature, predicted, actual, cprob=None
53
54 classifier.manualClassdb(num, title, feature, predicted, actual)
55
56 #def autoClassify(feed,classifier):
57 num=50
58 print '----- Begin automatic classification -----'
59 # Get feed entries and loop over them
60 f=feedparser.parse(feed)
61 for entry in f['entries'][50:100]:
62     num=num+1
63     # Print the contents of the entry
64     title=entry['title'].encode('utf-8').replace("'",'')
65     print 'Title:      '+ title
66     summary = splitRegexp.sub( "", entry[ "summary" ] )
67
68     print summary #entry['summary'].encode('utf-8')
69
70 # Combine all the text to create one item for the classifier
71 #fulltext='%s\n%s\n%s' % (entry['title'],entry['publisher'],entry['summary'])
72 fulltext='%s\n%s' % (entry['title'],entry['summary'])
73 fulltext=fulltext.replace("'",'')
74 # Print the best guess at the current category
75 predicted=str(classifier.classify(fulltext))
76 print 'Predicted: ', predicted
77
78 # Ask the user to specify the correct category
79 actual=raw_input('Enter actual category: ')
80 feature=raw_input('Enter string classifier: ')
81
82 #classifier.train(entry,cl)
83 # probability the item should be in this category
84 cp=round(classifier.cprob(feature, predicted),3)

```

```

85     fp=round(classifier.fisherprob(feature,predicted),3)
86     # print 'cprob: ', str(cp)
87     print 'fisherprob: ', str(fp)
88
89     # Save the trained classifications
90     # num, entry, feature, predicted, actual, cprob(feature, predicted)
91     classifier.autoClassdb(num, title, feature, predicted, actual, cp)
92     #return classifier
93
94 def entryfeatures(entry):
95     splitter=re.compile('\W*')
96     f={}
97
98     # Extract the title words and annotate
99     titlewords=[s.lower() for s in splitter.split(entry['title'])
100                 if len(s)>2 and len(s)<20]
101     for w in titlewords: f['Title:'+w]=1
102
103     # Extract the summary words
104     summarywords=[s.lower() for s in splitter.split(entry['summary'])
105                  if len(s)>2 and len(s)<20]
106
107     # Count uppercase words
108     uc=0
109     for i in range(len(summarywords)):
110         w=summarywords[i]
111         f[w]=1
112         if w.isupper(): uc+=1
113
114     # Get word pairs in summary as features
115     if i<len(summarywords)-1:
116         twowords=' '.join(summarywords[i:i+1])
117         f[twowords]=1
118
119     # Removed: Keep creator and publisher whole
120     #f['Publisher:'+entry['publisher']]=1
121
122     # UPPERCASE is a virtual word flagging too much shouting
123     if float(uc)/len(summarywords)>0.3: f['UPPERCASE']=1
124
125     return f

```

2.3.2 prog.py

```
1 import feedfilter
2
3 def main():
4     cl=docclass.fisherclassifier(docclass.getwords)
5     cl.setdb('allsports.db')
6     read('allsports.xml',cl)
7
8 if __name__ == "__main__":
9     main()
```

2.3.3 docclass.py

```
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     #print doc
9     ## Remove all the HTML tags
10    doc=re.compile(r'<[>]+>').sub('',doc)
11    # Split the words by non-alpha characters
12    words=[s.lower() for s in splitter.split(doc)
13           if len(s)>2 and len(s)<20]
14
15    # Return the unique set of words only
16    return dict([(w,1) for w in words])
17
18 class classifier:
19     def __init__(self,getfeatures,filename=None):
20         # Counts of feature/category combinations
21         self.fc={}
22         # Counts of documents in each category
23         self.cc={}
24         ## extract features for classification
25         self.getfeatures=getfeatures
26
27     def setdb(self,dbfile):
28         self.con=sqlite.connect(dbfile)
29         self.con.execute('create table if not exists rss(num, entry, feature, predicted,
30         actual, cprob)')
31         self.con.execute('create table if not exists fc(feature,category,count)')
32         self.con.execute('create table if not exists cc(category,count)')
33         # remove old data from previous sessions
34         # self.con.execute('delete from rss')
35         # self.con.execute('delete from fc')
36         # self.con.execute('delete from cc')
37
38     def manualClassdb(self,num, entry, feature, predicted, actual):
39         self.con.execute("insert into rss values ('%s','%s', '%s', '%s','%s', '%s')"%
40         % (num, entry, feature, predicted, actual, None))
41         self.con.commit()
42
43     def autoClassdb(self,num, entry, feature, predicted, actual, cp):
```

```

43     self.con.execute("insert into rss values ('%s','%s', '%s', '%s','%s', '%s')")
44                     % (num, entry, feature, predicted, actual, cp))
45     self.con.commit()
46     ## Increase the count of a feature/category pair
47     def incf(self,f,cat):
48         count=self.fcount(f,cat)
49         if count==0:
50             self.con.execute("insert into fc values ('%s','%s',1)"
51                             % (f,cat.lower()))
52         else:
53             self.con.execute(
54                 "update fc set count=%d where feature='%s' and category='%s'"
55                 % (count+1,f,cat.lower()))
56
57     ## The number of times a feature has appeared in a category
58     def fcount(self,f,cat):
59         res=self.con.execute(
60             'select count from fc where feature="%s" and category="%s"'
61             % (f,cat)).fetchone()
62         if res==None: return 0
63         else: return float(res[0])
64
65     ## Increase the count of a category
66     def incc(self,cat):
67         count=self.catcount(cat)
68         if count==0:
69             self.con.execute("insert into cc values ('%s',1)" % (cat.lower()))
70         else:
71             self.con.execute("update cc set count=%d where category='%s'"
72                             % (count+1,cat))
73
74     ## The number of items in a category
75     def catcount(self,cat):
76         res=self.con.execute('select count from cc where category="%s"'
77                             % (cat)).fetchone()
78         if res==None: return 0
79         else: return float(res[0])
80
81     ## The list of all categories
82     def categories(self):
83         cur=self.con.execute('select category from cc');
84         return [d[0] for d in cur]
85
86     ## The total number of items
87     def totalcount(self):
88         res=self.con.execute('select sum(count) from cc').fetchone();
89         if res==None: return 0
90         return res[0]
91
92
93     ## The train method takes an item(document) and a classification.
94     ## It uses the getfeatures function to the break the item into its
95     ## separate features. It then calls incf to increase the counts for
96     ## this classification for every feature. Finally, it increases
97     ## the total count for this classification.
98     def train(self,item,cat):
99         features=self.getfeatures(item)

```

```

100     # Increment the count for every feature with this category
101     for f in features:
102         self.incf(f,cat)
103
104     # Increment the count for this category
105     self.incc(cat)
106     self.con.commit()
107
108     ## Probability is a number between 0 and 1, indicating
109     ## the likelihood of an event. You calculate the probability of
110     ## a word in a particular category by dividing the number of
111     ## times the word appears in a document in that category
112     ## by the total number of documents in the category.
113     def fprob(self,f,cat):
114         if self.catcount(cat)==0: return 0
115
116         # The total number of times this feature appeared in this
117         # category divided by the total number of items in this category
118         return self.fcount(f,cat)/self.catcount(cat)
119
120     def weightedprob(self,f,cat,prf,weight=1.0,ap=0.5):
121         # Calculate current probability
122         basicprob=prf(f,cat)
123
124         # Count the number of times this feature has appeared in
125         # all categories
126         totals=sum([self.fcount(f,c) for c in self.categories()])
127
128         # Calculate the weighted average
129         bp=((weight*ap)+(totals*basicprob))/(weight+totals)
130         return bp
131
132
133
134
135     class naivebayes(classifier):
136
137         def __init__(self,getfeatures):
138             classifier.__init__(self,getfeatures)
139             self.thresholds={}
140
141         def docprob(self,item,cat):
142             features=self.getfeatures(item)
143
144             # Multiply the probabilities of all the features together
145             p=1
146             for f in features: p*=self.weightedprob(f,cat,self.fprob)
147             return p
148
149         def prob(self,item,cat):
150             catprob=self.catcount(cat)/self.totalcount()
151             docprob=self.docprob(item,cat)
152             return docprob*catprob
153
154         def setthreshold(self,cat,t):
155             self.thresholds[cat]=t
156

```

```

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

```



```

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

```

3 Question 3:

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

3.1 Approach

Precision, Recall and F measure are based on TP(True positive), True Negative(TN), False Positive(FP), False Negative(FN)

- Precision is the fraction of retrieved instances that are relevant
- Precision= $TP / (TP + FP)$
- Recall is the fraction of relevant instances that are retrieved
- Recall= $TP / (TP + FN)$
- F-Measure is a measure of a test's accuracy. It considers both the precision and the recall
- $F1 = 2TP / (2TP + FP + FN)$

3.2 Tables

Category	TP	TN	FP	FN
cricket	19	6	1	24
olympics	1	0	4	45
football	16	4	0	30
tennis	0	0	3	47
hockey	0	0	1	49
others	1	1	3	45

Table 3: TP TN FP FN values of different categories

Category	Precision	Recall	F1
cricket	0.95	0.44	0.6031
olympics	0.2	0.02173	0.03921
football	1	0.3478	0.5161
tennis	0	0	0
hockey	0	0	0
others	0.25	0.02173	0.04

Table 4: Precision Recall and F1 values for each category

References

- [1] F-Measure wiki: https://en.wikipedia.org/wiki/F1_score.
- [2] Precision Recall Wiki: https://en.wikipedia.org/wiki/Precision_and_recall.
- [3] Programming Collective Intelligence : <https://github.com/manojchandrak/Programming-Collective-Intelligence>.

[]