Assignment 9
CS 532: Introduction to Web Science Spring 2018 Chandrasekhar Reddy Muthyala

1

Question

- 1. Using the data from A8:
- Consider each row in the blog-term matrix as a 1000 dimension vector, corresponding to a blog.
- From chapter 8, replace numpredict.euclidean() with cosine as the distance metric. In other words, you'll be computing the cosine between vectors of 1000 dimensions.
- Use knnestimate() to compute the nearest neighbors for both:

```
http://f-measure.blogspot.com/
http://ws-dl.blogspot.com/
```

for $k=\{1,2,5,10,20\}$.

Answer

To solve this question I used the blog data I previously received in assignment 7 blogdata.txt and the code provided by the Programming Collective Intelligence book [2].

Euclidean distance is probably the one that you are most familiar with; it is essentially the magnitude of the vector obtained by subtracting the training data point from the point to be classified. Another common metric is Cosine similarity. Rather than calculating a magnitude, Cosine similarity instead uses the difference in direction between two vectors [4].

$$similarity = COS(\Theta) = \frac{A.B}{||A||||B||}$$

To calculate the cosine distance metric **numpredict.py** had to be modified to accompany this new function. The methods that changed in this file were the *knnestimate* and *getdistances* functions. In the *getdistances* function I simply swapped out the euclidean function for cosine as shown in Listing 1. In the *knnestimate* function and return a list of sorted distances in descending order.

I then went filtered the F-measure and Web Science research group's blogs and create separate vectors with their values. These vectors were used for the cosine measurement against the vector with all blogs.

```
1
2
   def cosine(v1, v2):
3
        cosine function for 1000 blog matrix. zip v1, v2
4
5
        for array of tuples
6
7
        sumab = sum([a * b for a, b in zip(v1, v2)])
8
        suma = sum([a * a for a in v1])
9
        sumb = sum([b * b for b in v2])
10
        return sumab / math.sqrt(suma * sumb)
11
12
13
   def getdistances (data, vec1):
14
        distancelist = []
15
        # Loop over every item in the dataset
16
17
        for i, vals in enumerate (data):
            vec2 = vals
18
19
```

```
20
            # Add the distance and the index
21
            distancelist.append((cosine(vec1, vec2), i))
22
23
        # Sort by distance
24
        distancelist.sort(reverse=True)
25
        return distancelist
26
27
28
    def knnestimate(data, vec1, k=5):
29
        # Get sorted distances
        distances = getdistances (data, vec1)
30
31
        return distances
32
33
34
    def inverseweight (dist, num=1.0, const=0.1):
        \texttt{return num / (dist + const)}
35
36
37
38
    def subtractweight (dist, const=1.0):
39
        if dist > const:
40
            return 0
```

Listing 1: Python script with included cosine function and knnestimate changes

```
1
   from numpredict import *
3
4
   def\ estimate (\,vectorValues\,,\ vectorBlog1\,,\ vectorBlog2\,):
5
6
       nn = knnestimate(allBlogVector.values(), vectorBlog1)
7
       8
       kvals = [1, 2, 5, 10, 20]
9
       for k in kvals:
10
          print('k = ', k)
11
           for j in range(k):
              print('%s' % (list(allBlogVector.keys())[nn[j][1]]))
12
          print("----")
13
14
15
       print ("*************tWeb Science and Digital Libraries
          16
       nn = knnestimate(allBlogVector.values(), vectorBlog2)
17
       for k in kvals:
           print('k = ', k)
18
           for j in range(k):
19
           print('%s' % (list(allBlogVector.keys())[nn[j][1]]))
print("----")
20
21
22
23
```

```
24
25
   def formatData():
        givenBlog_1 = 'F-Measure'
26
27
        givenBlog_2 = 'Web Science and Digital Libraries Research
28
        allBlogVector = \{\}
        blog1List = []
29
30
        blog2List = []
31
        with open("blogdata.txt", 'r') as f:
32
            allLines = f.readlines()
            for i, line in enumerate(allLines):
33
34
                if i == 0:
35
                    # skip header
36
                    continue
37
                blogMatrix = line.strip().split('\t')
                if blogMatrix[0] = givenBlog_1:
38
39
                    for i in range(1, len(blogMatrix)):
40
                         blog1List.append(float(blogMatrix[i]))
41
                elif blogMatrix[0] = givenBlog_2:
                    for i in range(1, len(blogMatrix)):
42
43
                         blog2List.append(float(blogMatrix[i]))
                else:
44
                    allBlogVector[blogMatrix[0]] = []
45
46
                    for i in range(1, len(blogMatrix)):
47
                         allBlogVector[blogMatrix[0]].append(float(
                            blogMatrix[i]))
        return allBlogVector, blog1List, blog2List
48
49
50
   if -name_{-} = "-main_{-}":
51
        allBlogVector, vectorBlog1, vectorBlog2 = formatData()
52
53
        estimate(allBlogVector.values(), vectorBlog1, vectorBlog2)
```

Listing 2: Python script to find KNN values

OutPut:

```
11 \mid ('k = ', 5)
12
   DaveCromwell Writes
13
   The Jeopardy of Contentment
   Pithy Title Here
14
   KiDCHAIR
15
16
   Alex Denney
17
   ('k = ', 10)
18
   DaveCromwell Writes
19
20
   The Jeopardy of Contentment
21
   Pithy Title Here
22
   KiDCHAIR
23
    Alex Denney
24
    Skiptrack music
25
   New Amusements
26
   My Name Is Blue Canary
27
   The Slow Music Movement
28
29
30
  ('k = ', 20)
  DaveCromwell Writes
31
32
   The Jeopardy of Contentment
33
   Pithy Title Here
34
   KiDCHAIR
35
   Alex Denney
36
   Skiptrack music
37
   New Amusements
38
   My Name Is Blue Canary
39
   The Slow Music Movement
40
   Diagnosis: No Radio
41
   The Power of Independent Trucking
42
   Did Not Chart
44
   Captain Panda's Local & Independent Music Showcase
45
   PSI LAB
46
   unter diesem gesichtspunkt
47
   Lyrically Speaking
48
   macthemost
49
   Notes from a Genius
50
   She May Be Naked
51
52
    ******
                    Web Science and Digital Libraries Research Group
    ('k = ', 1)
53
54
   unter diesem gesichtspunkt
55
56
   ('k = ', 2)
   unter diesem gesichtspunkt
57
58 | Myopiamuse
```

```
59
    ('k = ', 5)
60
    unter diesem gesichtspunkt
61
62
    Myopiamuse
    She May Be Naked
63
    Alex Denney
65
    DaveCromwell Writes
66
67
   ('k = ', 10)
68
    unter diesem gesichtspunkt
    Myopiamuse
69
70 She May Be Naked
71
    Alex Denney
72
    DaveCromwell Writes
73
    Pithy Title Here
74
    The Power of Independent Trucking
    The Professional Daydreamer
75
76
    Diagnosis: No Radio
77
    My Name Is Blue Canary
78
79 ('k = ', 20)
80 unter diesem gesichtspunkt
81 Myopiamuse
82 | She May Be Naked
83
    Alex Denney
   DaveCromwell Writes
85
    Pithy Title Here
86
    The Power of Independent Trucking
87
    The Professional Daydreamer
    Diagnosis: No Radio
88
    My Name Is Blue Canary
89
    Mile In Mine
    Morgan's Blog
92
    Media Coursework
93
94
    macthemost
95
   New Amusements
96
    ELLIA TOWNSEND A2
97
    hello my name is justin.
98
    PSI LAB
99
    Nothing But Ordinary Glances At Extraordinary Things
100
```

Listing 3: Output of the knnestimate for two given blogs

References

- [1] https://github.com/cmuth001/anwala.github.io/tree/master/Assignments/A7.
- [2] Segaran, Toby. "Programming Collective Intelligence". O' Reilly, 2007. Web. 6 April 2017. http://shop.oreilly.com/product/9780596529321.do.
- [3] https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm
- [4] (https://towardsdatascience.com/introduction-to-k-nearest-neighbors-3b534bb11d26)