# CS-432/532 Introduction to Web Science: Assignment #10: kNN and SVM

Due on Saturday, April 30, 2016

Dr. Michael L. Nelson

Plinio Vargas pvargas@cs.odu.edu

# Contents

	em 1         Approach          Solution	1 1 4
List	of Figures	
Listi	ings	
1	Cosine Calculator: CosineDistance.py	1
2	Main Module: CosineDistance.py	
3	Getdistances Module: CosineDistance.py	2
4	Knnestimate Module: CosineDistance.py	3
List	of Tables	
1	F-Measure with K=1	9
2	F-Measure with $K=2$	9
3	F-Measure with K=5	9
4	F-Measure with K=10	10
5	F-Measure with K=20 $\dots$	10
6	dl.blogspot.com with K=1	11
7	dl.blogspot.com with K=2	11
8	dl.blogspot.com with K=5	11
9	dl.blogspot.com with K=10 $\dots$	11
10	dl.blogspot.com with K=20 $\dots$	12

# Problem 1

Using the data from A8:

- Consider each row in the blog-term matrix as a 500 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 500 dimensions.
- Use knnestimate() to compute the nearest neighbors for both:

```
\label{eq:http://f-measure.blogspot.com/http://ws-dl.blogspot.com/} $$ for $k=\{1,2,5,10,20\}. $$
```

## 1.1 Approach

In order to complete this problem, three modules from [1] were taken and modified: **knnestimate**, **get-distances** and *euclidean*, the latest changed to **cosine\_distance**. The heart of this approach resides in the **cosine\_distance** module (see Listing 1), which takes as a parameter two 500-dimensional vectors.

Listing 1: Cosine Calculator: CosineDistance.py

```
def cosine_distance(v1, v2):
    d = 0.0
    for i in range(len(v1)):
        d += (v1[i] - v2[i])
    return 1 - math.cos(d)
```

Since  $\cos \theta$  for any point is equal to 1 when x and y are equal in our plotting coordinate, if the cosine distance calculation between vector  $v_1$  and  $v_2$  is equal to zero, then  $v_1 = v_2$ . The further from 0 the calculation between  $v_1$  and  $v_2$  is, the greater the difference is between them.

Using the same approach as in [1], the values are loaded into a dictionary object (data). Term values are extracted from the first line of the input file (see Listing 2 line 21).

Listing 2: Main Module: CosineDistance.py

```
for record in file:
18
               record = record.split()
               if k < 1:
20
                   terms = record[-500:]
               else:
                    print(k, len(record) - m, record[0:len(record) - m])
23
                   title.append(record[0:len(record) - m])
24
                   data.append({'input': tuple([int(x) for x in record[-500:]])})
25
               k += 1
      print(terms)
27
28
      for k in range(100):
29
           print('(%s, %d)' % (' '.join(title[k]), k), end=',')
30
      print()
31
32
      k_{values} = [1, 2, 5, 10, 20]
      print('Closest Blog to %s' % ' '.join(title[76]))
34
      for k in k_values:
35
           print('\nk=%d' % k)
36
           knnestimate(data, data[76]['input'], k)
37
38
      print()
39
      print('Closest Blog to %s' % ' '.join(title[68]))
      for k in k_values:
41
42
           print('\nk=%d' % k)
           knnestimate(data, data[68]['input'], k)
44
48
      return
```

The blog's title (line 24) is obtained by removing the 500-dimensional vector values from the file entry row. To obtain the row index from our two studied blogs, we visually inspected the entry file (blog-data.tx). Index 76 corresponded to http://f-measure.blogspot.com/; index 68 corresponded to http://ws-dl.blogspot.com/

In order to get the cosine-distance between two blogs, we passed the vector values to function *getdistances*. As the textbook explains, it "calls the distance function on the vector given against every other vector in the dataset and puts them in a big list. The list is sorted so that the closest item is on the top".

Listing 3: Getdistances Module: CosineDistance.py

```
def getdistances(data, vec1):
    distancelist = []

# Loop over every item in the dataset
for i in range(len(data)):
    vec2 = data[i]['input']

# Add the distance and the index
    distancelist.append((cosine_distance(vec1, vec2),i))

# Sort by distance
distancelist.sort()
```

```
13
14 return distancelist
```

To display k-blogs closest to blog-x, an iteration is performed to call function knnestimate. This function, a modification from [1], simply returns the k-elements from the sorted list provided by function get distances between blog-x and the other 99 blogs.

Listing 4: Knnestimate Module: CosineDistance.py

```
def knnestimate(data, vec1, k=5):
    # Get sorted distances
    dlist = getdistances(data, vec1)
    avg = 0.0

print(dlist)
    # Take the average of the top k results
    for i in range(1, k + 1):
        idx = dlist[i][1]
        print(' '.join(title[idx]))

return avg
```

#### 1.2 Solution

Since we clustered our 100 blogs in assignment 8, we can use the previous work as a comparison. Tables 1 through 10 contain our two blog similarities among the selected 100 blogs, using cosine-distance calculator.

According to Table 1, the most similar blog to **F-Measure** is http://didnotchart.blogspot.com/. This blog passed the eye test since it is also a blog related to music. Looking at the dendrogram from assignment 8, "Did Not Chart" is clustered close to **F-Measure**, but is not the closest. **F-Measure** was clustered closer to http://ihatethe90s.blogspot.com/. The last is also a blog related to music, so it also passed the eye test.

Then, which blog is the most similar to **F-Measure**? http://ihatethe90s.blogspot.com/ or http://didnotchart.blogspot.com/? To answer this question we need to dive into the raw data, the DNA of our research:

#### DidNotChart

 $2011058000030020120100000000132022016002000110710140230401110000213323100000002822\\ 1102104201470307210000271830001010130110000007300240105000000011301301112000042010\\ 0061105200020020011006000050060012000101002110210320314001000004000110100002206011\\ 0090010110300003131003420100000400620901001000008201031120110100014103030105538042\\ 0301011000000110301120411100000023000020003000102000001360000110300410041200317001\\ 10105120020018000010511501103103100010010000740071250010213002004030100101050110\\ 110330000000110000$ 

#### F-Measure

 $200100038011111136210004000010001013000000109010200210000010004000012032000200251\\0000010041090502001010081260312211022400100302040005030110005101020020200212124040\\21112144800002000012208000110010011110000101020210183022020103001101234020021821031\\0120051016004100020030201000001010300011240012000110161011002200200114600000212222\\0221002031422000150202025030330110000505100001110020000100100015000602000014140101\\0130600118007100010001600110001080110200000012142000005001025003016010101140030101\\20500130200212330060$ 

#### IHateThe90s

 $5410143611006300011153261001450341014154144103110204150232112250740402910902119131\\1274001966276235740638572282415100014910212011827014312221667521601610020114504224\\3503713405016477401112118800214131713025333403030220527161217310458621330327411261\\2339112264622111360142121313225128312673441211131110146111120093036121633230260135\\5224414116022175331910320211641030125491745601965317211133914122066001526116214130\\3116232430346315842341237301310760156050201822631080330161565567576122811132101343\\0434176111324110139300122076451710610601512052912135803628211270000001121511821215\\216531111147034$ 

To understand the raw data above, we need to inspect the 500 matrix file: blogdata.txt. The DNA thickness difference between I-Hate-The-90s and F-Measure is not due to extra terms, but rather many term frequency in I-Hate-The-90s. The differences are in the double digits, while in F-Measure and Did-Not-Chart, the differences are in the single digits. Then, although the blogs could be talking about the same thing, their frequencies are so far apart that they make them less similar.

As we expand the number of blog k-neighbors for K=20, an examination to the most distant neighbors "funky little demons" revealed this blog is not exactly about music, but it has lyrics which are related to

music. Looking at dendrogram from assignment 8, "funky little demons" was not clustered near **F-Measure**, and the deepest element in the cluster is http://kidchair.blogspot.com/, which is about songs (music).

#### F-Measure

 $200100038011111136210004000010001013000000109010200210000010004000012032000200251\\0000010041090502001010081260312211022400100302040005030110005101020020200212124040\\21112144800002000012208000110010011110000101020210183022020103001101234020021821031\\012005101600410002003020100000101300011240012000110161011002200200114600000212222\\0221002031422000150202025030330110000505100001110020000100100015000602000014140101\\0130600118007100010001600110001080110200000012142000005001025003016010101140030101\\20500130200212330060$ 

#### funkylittledemons

#### **KiDCHAIR**

The answer to the question which one is more similar to **F-Measure** between the last two blogs is more difficult to answer from a 40K-feet-view looking at their DNAs. The frequency of 500 terms in "funky little demons" is very low; many elements in the vector have zero values. We could be very quick to say that http://kidchair.blogspot.com/ is more similar to **F-Measure** than "funky little demons", but if we were going element by element in the vector, we will notice the combined difference of the last blog is smaller.

A surprising result was the comparison with http://ws-dl.blogspot.com/. The closest blog in similarity was http://lostintheshuffle899.blogspot.com/ - Lost in the Shuffle (see Table 6); however in the dendrogram from assignment 8, it was clustered closer to **SpinitronBlog**.

Using similar comparison to their DNAs, we can notice **SpinitronBlog** does not have too many terms while **LostintheShuffle** does, making it closer by default.

#### WebScienceandDigitalLibrariesResearchGroup

 $3800605137033020145201800443102040118128020060021020010112007011207710014600000031\\7131702030031167060300019306001230101301900040113101800001208002000010100211060413\\3075030204511041112157802890560271238200024611100723112092600200015630000007370010\\6843114500301131012124200206105100050201310880020190100050002000860133022000000031\\2615000101340100010200231360210011150140004014172012242045011200011501183800022121$ 

 $0254012010003512030003000221310206097212151313220710130201101111001792541113150270\\0022001000006070201243010011000010019180101640004105100103130213390100050212$ 

#### LostintheShuffle

#### SpinitronBlog

### Blog-Terms:

groups-ashtray-stretching-sum-indians-adventure-rightly-cell-reviewing-tstyle-reversestatements-chilling-existent-sleepy-fortunately-typically-settle-exit-lone-behavecompositions-bastards-canned-pagination-jacob-reflective-legit-eleven-exclusivelyspending-rival-wings-license-norm-encourage-haircut-emerge-actors-beg-honored-closetoutfits-pr-collector-whore-teenagers-leisure-mob-ali-marry-hopeful-jobs-knocking-westernspleased-females-elliot-someday-cutting-ken-historic-harris-meditation-insane-chant-mythchallenge-passive-longing-curse-lenny-jewelry-ducks-supreme-joan-graduation-nigelrecordings-pedal-security-tale-butter-redhead-spice-oklahoma-html-ringing-combine-armysentiments-funk-cuff-violence-orbit-producer-ariel-switched-snakes-careless-equation-arrivinglarger-syd-trivia-sharp-quit-critic-benefits-norton-colbert-glen-industry-clue-appearedrivers-shaking-recipe-expanded-undoubtedly-endings-dancer-sailing-viewed-jurado-drove-kilsheen-cassettes-steam-twigs-rockabilly-pummeling-reminder-sensibility-practically-lame-granderesulting-chatter-afghan-remembering-breathtaking-carl-alps-happiness-buying-carrie-careersdual-bucks-origins-quintessential-infinite-andrews-sunny-tower-chapters-ipod-immortal-mahistorically-cinema-escaped-abuse-powered-eric-argue-lip-approaches-schools-cramps-ash-bloggedgentleman-caring-downs-cos-appropriate-attempting-accepting-mo-eliminate-polar-cornell-flashchristian-leaning-electricity-strongly-grief-polite-responsible-rebellion-fake-rides-implies-hamelvins-remove-plethora-gained-bundle-mc-crashed-astronaut-celebrate-fists-patient-knowledgepulled-wives-symphony-laundry-instinct-spiritualized-mister-allow-landed-minority-futuristiccatches-surprises-claimed-index-railroad-fictional-luke-importance-aquarium-daddy-neighborsfourteen-dime-kim-sports-delay-loser-slower-kramer-gutter-grinding-pro-gods-vocabulary-moviesgoal-mondays-birth-mould-partnership-skinny-upcoming-grandfather-delivering-groundbreakingjacuzzi-skills-announcement-insanely-sprawling-lucas-discography-served-humming-ashley-idiotscrows-truckers-subsequently-amused-effortlessly-haze-additions-flourishes-mccarthy-heathexecution-formed-tons-analysis-winding-percentage-evokes-donuts-impress-outlet-zone-survivingjulian-complain-wheels-palace-edward-shaped-ratio-tones-ry-cred-trapped-puppets-simpson-insertconsume-zeppelin-virtue-eps-nostalgic-alaska-bend-swallowed-accidentally-hallelujah-russiadirected-misery-greg-bizarre-dean-exhausting-eccentric-formerly-clicks-accurate-crashes-modrarity-goodies-kenny-shannon-pros-categories-majority-audible-cosmic-euro-proven-product-remedyelliott-problematic-amongst-slate-blew-expert-tribe-circus-vox-iris-ny-ve-crushed-cartoonparties-university-hyper-sidewalk-jan-trials-promised-wisconsin-grape-victims-matched-graphicbooker-umbrella-eternal-junior-swamp-oil-drake-originality-apartment-potter-marianne-explosivedum-lily-chromatics-magnet-disgusting-pockets-stanley-assured-oldham-illustrated-legend-adslogo-continent-noir-respond-tribal-sites-barnes-delivery-impression-health-info-chairs-losstraw-directors-desk-simone-thin-cheaper-thrive-allen-maiden-elite-ab-combined-venue-jonnynothings-emergency-axis-studied-necessity-bush-confirm-newest-walker-shades-windows-realistichandy-mystic-correct-palmer-eden-sans-observed-coachella-gimmick-fashioned-traded-titus-elliscouch-bean-bursts-labor-quietly-featured-tease-species-louise-guessed-reviewer-experimental-msoverhead-breakdown-drifting-amy-finest-northwest-babes-lonesome-owls-korea-bow-swordsconfidence-suzanne-observations-trumpet-pile-hating-reflects-websites-horribly-cheating-sufjanintensity-advantage-engage-zombie-laced-moore-wears-tremendous-winners-songwriters-concisetemptation-entering-jefferson-casey-veteran-coincidentally-distorted-riding-worlds-germanslang-manning

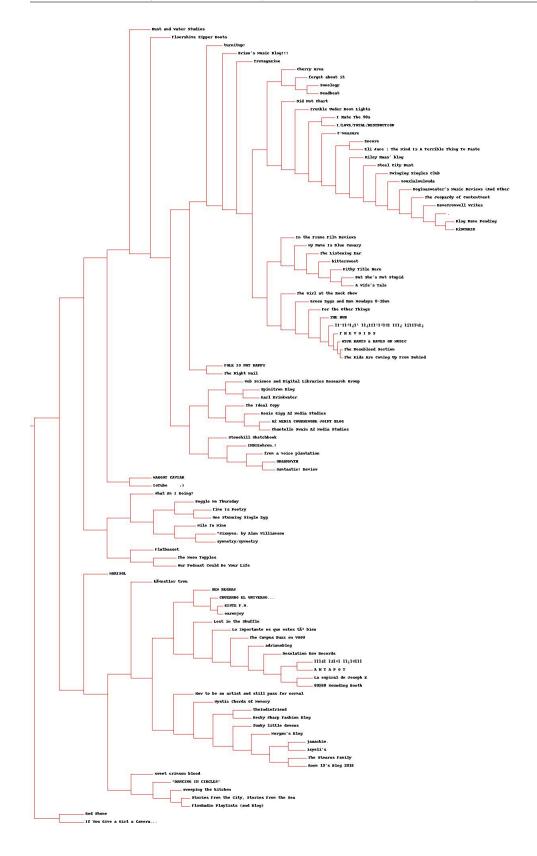


Table 1: F-Measure with K=1

Blog Title
Did Not Chart

Closest K-Blogs to F-Measure for K=1. Values obtained from CosineDistance.py

Table 2: F-Measure with K=2

Blog Title
Did Not Chart
THIS CHARMING YAN

Closest K-Blogs to F-Measure for K=2. Values obtained from CosineDistance.py

Table 3: F-Measure with K=5

Blog Title
Did Not Chart
THIS CHARMING YAN
Encore
isyeli's
If You Give a Girl a Camera...

Closest K-Blogs to F-Measure for K=5. Values obtained from CosineDistance.py

Table 4: F-Measure with K=10

Blog Title

Did Not Chart

THIS CHARMING YAN

Encore

isyeli's

If You Give a Girl a Camera...

mattgarman

Rants from the Pants

60@60 Sounding Booth

turnitup!

Doginasweater's Music Reviews (And Other ..)

Closest K-Blogs to F-Measure for K=10. Values obtained from CosineDistance.py

Table 5: F-Measure with K=20

Blog Title

Did Not Chart

THIS CHARMING YAN

Encore

isyeli's

If You Give a Girl a Camera...

mattgarman

Rants from the Pants

60@60 Sounding Booth

turnitup!

Doginasweater's Music Reviews (And Other ..)

"DANCING IN CIRCLES"

I/LOVE/TOTAL/DESTRUCTION

Lost in the Shuffle

this time tomorrow

Web Science and Digital Libraries Research Group

FlowRadio Playlists (and Blog)

.

Stories From the City, Stories From the Sea

The Jeopardy of Contentment

funky little demons

Closest K-Blogs to F-Measure for K=20. Values obtained from CosineDistance.py

Table 6: dl.blogspot.com with K=1

Blog Title
Lost in the Shuffle

Closest K-Blogs to dl.blogspot.com for K=1. Values obtained from CosineDistance.py

Table 7: dl.blogspot.com with K=2

Blog Title
Lost in the Shuffle
"DANCING IN CIRCLES"

Closest K-Blogs to dl.blogspot.com for K=1. Values obtained from CosineDistance.py

Table 8: dl.blogspot.com with K=5

Blog Title

Lost in the Shuffle

"DANCING IN CIRCLES"

Stories From the City, Stories From the Sea Doginasweater's Music Reviews (And Other Horse ...) turnitup!

Closest K-Blogs to dl.blogspot.com for K=5. Values obtained from CosineDistance.py

Table 9: dl.blogspot.com with K=10

Blog Title

Lost in the Shuffle

"DANCING IN CIRCLES"

Stories From the City, Stories From the Sea

Doginasweater's Music Reviews (And Other Horse ...)

turnitup!

The Jeopardy of Contentment

If You Give a Girl a Camera...

mattgarman

Pop Tones

isyeli's

Closest K-Blogs to dl.blogspot.com for K=10. Values obtained from CosineDistance.py

Table 10: dl.blogspot.com with K=20

## Blog Title

Lost in the Shuffle

"DANCING IN CIRCLES"

Stories From the City, Stories From the Sea

Doginasweater's Music Reviews (And Other Horse ...)

turnitup!

The Jeopardy of Contentment

If You Give a Girl a Camera...

mattgarman

Pop Tones

isyeli's

Encore

THIS CHARMING YAN

Samtastic! Review

Did Not Chart

F-Measure

Room 19's Blog 2016

Rants from the Pants

60@60 Sounding Booth

Cherry Area

I/LOVE/TOTAL/DESTRUCTION

 $\label{eq:closest} \text{Closest K-Blogs to dl.blogspot.com for K=20.} \text{Values obtained from CosineDistance.py}$ 

# References

- [1] Segarn, Toby. Programming Collective Intelligence. Building Smart Web 2.0 Application. (pp 29-53). Sebastopol, CA: O'Reilly Media.
- [2] Text Mining, Analytics & More. (n.d.) Retrieved April 21, 2016, from http://www.text-analytics101.com/2014/10/computing-precision-and-recall-for.html