



DSC478: Programming Machine Learning Applications

Roselyne Tchoua

rtchoua@depaul.edu

School of Computing, CDM, DePaul University

General Notes

- Submit one notebook per question for clarity
- Zip all files and submit a single file

General Notes

- Remember to use markdown section for any observations. Comments can be used for short code comments but make sure to annotate any plot/table and reply separately to any question asking for your comments.
 - For example, you are asked to compare results with Euclidean and Cosine similarities, or results with and without TFIDF weights, use markdown cells to make observations

Compare/contrast/discuss

- You have definition for bias and variance, you should be able to not only define these concepts but observe them in your results. Especially variance = difference in accuracy with different training sets.
- Your overall observations can be a combination of accuracy observations and what you answered above (bias-variance) with potential suggestions of what you would do next.

Data Format

Even though you would expect a full frame of training and testing, the data for the news dataset is already separated in training vs. testing data and training vs. testing labels. It is also not in the right format for our typical matrix, so you need to transpose it.

In [3]:

```
1 #Load data
2 train=pd.read_table('Data/trainMatrixModified.txt', header= None)
3 train_lab=pd.read_table('Data/trainClasses.txt', header=None, index_col=0)
4 test=pd.read_table('Data/testMatrixModified.txt', header=None)
5 test_lab=pd.read_table('Data/testClasses.txt', header=None, index_col=0)
6 terms = pd.read_table('Data/modifiedterms.txt', header=None)
```

In [4]:

```
1 train.head()
```

Out[4]:

	0	1	2	3	4	5	6	7	8	9	...	790	791	792	793	794	795	796	797	798	799
0	2.0	0.0	0.0	2.0	2.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
1	2.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	3.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2	2.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
3	1.0	1.0	1.0	1.0	1.0	1.0	2.0	1.0	1.0	1.0	...	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
4	8.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	2.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0

5 rows × 800 columns

In [5]:

```
1 train.shape
```

Out[5]: (5500, 800)

In [6]:

```
1 train_lab.shape
```

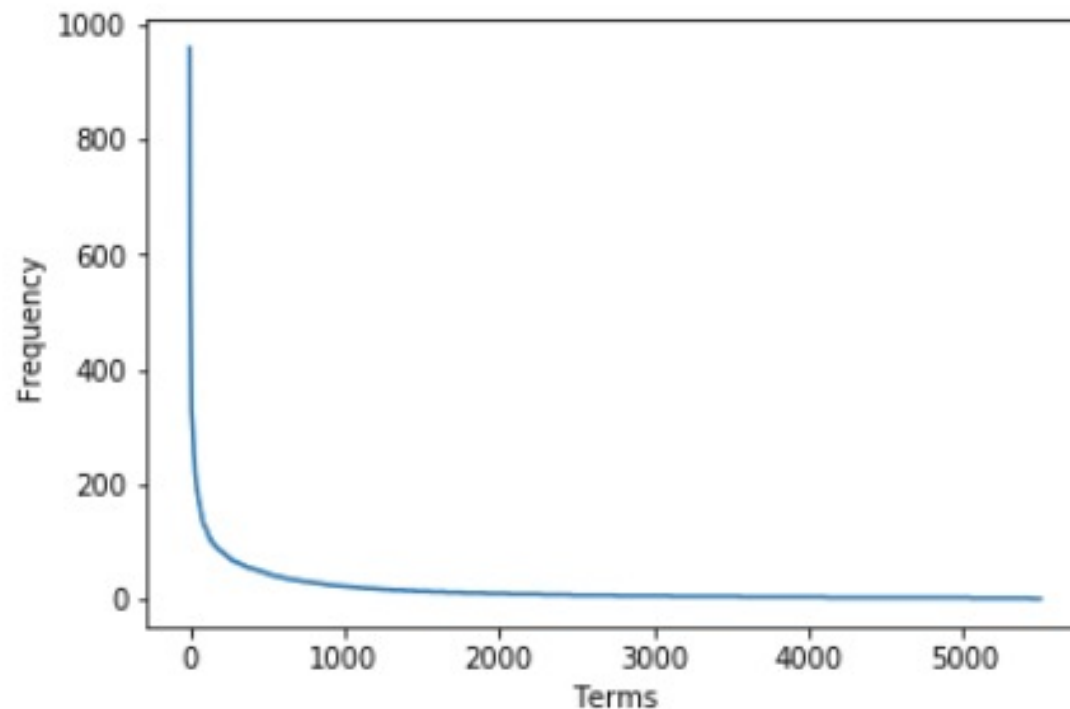
Out[6]: (800, 1)

In [7]:

```
1 #converting from TD to DT matrix
2 DT_train = train.T
3 DT_test = test.T
```

Zipf Distribution in Training Data

```
In [7]: plt.plot(sorted(termFreqs, reverse=True));  
plt.xlabel('Terms')  
plt.ylabel('Frequency')  
plt.show()
```



```
train_labels=pd.read_table('Data/trainClasses.txt', header=None, index_col=0)
```

```
In [97]: np.array(train_labels)
```

```
Out[97]: array([[0],  
                [1],  
                [0],  
                [1],  
                [0],  
                [1],  
                [1],  
                [1],  
                [1],  
                [1],  
                [1],  
                [1],  
                [0],  
                [0],  
                [1],  
                [1],  
                [1],  
                [0],  
                [0],  
                [1],  
                [0],
```



```
train_labels=pd.read_table('Data/trainClasses.txt', header=None, index_col=0)
```

```
In [98]: np.array(train_labels).flatten()
```

```
Out[98]: array([0, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 1,
0, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0,
0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 1, 1,
1, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0,
1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1,
1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1,
0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0,
1, 0, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 0, 1, 0,
1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0,
1, 0, 1, 1, 1, 1, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1,
0, 0, 1, 1, 1, 1, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1,
1, 1, 0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 1,
1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0,
0, 1, 1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1,
1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0,
0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0,
1, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0,
1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0,
0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1, 0, 0, 0, 1,
1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1,
0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1,
1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1,
0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1,
1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0,
0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0,
1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0,
0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0,
0, 0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 0, 1,
0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 0, 1, 1,
1, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 1, 1,
1, 0, 1, 0, 1, 1, 1, 1, 1], dtype=int64)
```

See the Class Notebook: [TF*IDF and Document Categorization](#)

```
1 def knn_search(x, D, K, measure):
2     """ find K nearest neighbors of an instance x among the instances in D """
3     if measure == 0:
4         # euclidean distances from the other points
5         dists = np.sqrt(((D - x)**2).sum(axis=1))
6     elif measure == 1:
7         # first find the vector norm for each instance in D, as well as for vector x
8         D_norm = np.array([np.linalg.norm(D[i]) for i in range(len(D))])
9         x_norm = np.linalg.norm(x)
10        # Compute Cosine: dot product of x and each instance in D divided
11        # by the product of the two norms
12        sims = np.dot(D,x)/(D_norm * x_norm)
13        # The distance measure will be the inverse of Cosine similarity
14        dists = 1 - sims
15        idx = np.argsort(dists) # sorting
16        # return the indexes of K nearest neighbors
17        return idx[:K], dists
18
19 def knn_classify(x, D, K, labels, measure):
20     from collections import Counter
21     neigh_idx, distances = knn_search(x, D, K, measure)
22     neigh_labels = labels[neigh_idx]
23     count = Counter(neigh_labels)
24     # print("Labels for top", K, "neighbors: ", count)
25     predicted_label = count.most_common(1)[0][0]
26     return neigh_idx, predicted_label
27
```

```
1 #preparation: transform data to matrix
2 DTM_train = np.array(DT_train)
3 DTM_test = np.array(DT_test)
4 train_lab_array = np.array(train_lab).flatten()
5 test_lab_array = np.array(test_lab).flatten()
```



```
1 # knn_function test of measure 1
2 top_K_neighbors, predicted_class = knn_classify(DTM_test[12], DTM_train, 10, train_lab_array, 1)
3 print ("The indices of K Nearest Neighbors are: {} \nThe predicted class is: {}".format(top_K_neighbors, predicted_class))
```

The indices of K Nearest Neighbors are: [334 177 233 715 732 510 221 178 94 254]
The predicted class is: 1

```
1 # knn_function test of measure 0
2 top_K_neighbors, predicted_class = knn_classify(DTM_test[12], DTM_train, 10, train_lab_array, 0)
3 print ("The indices of K Nearest Neighbors are: {} \nThe predicted class is: {}".format(top_K_neighbors, predicted_class))
```

The indices of K Nearest Neighbors are: [798 757 119 398 38 224 711 615 751 69]
The predicted class is: 0

Knn Classification & Evaluation

```
knn_search(x, DT_train, K, measure):
```

```
knn_classify(x, DT_train, K, train_labels, measure):
```

```
knn_evaluate(DT_test, test_labels, DT_train, train_labels, K, measure)
```

Note: **knn_evaluate** iterates through rows in DT_test matrix of test instances, and for each Instance, calls **knn_classify** to predict the label. It compares the predicted label for each test instance to the actual labels and then returns classification accuracy.

Knn Classification & Evaluation

```
1 Euclid=[]  
2 Cosine=[]
```

```
1 %%time  
2 for K in range(5, 100, 5):  
3     Euclid.append(knn_evaluate(DTM_test, test_lab_array, DTM_train, train_lab_array, K, 0))
```

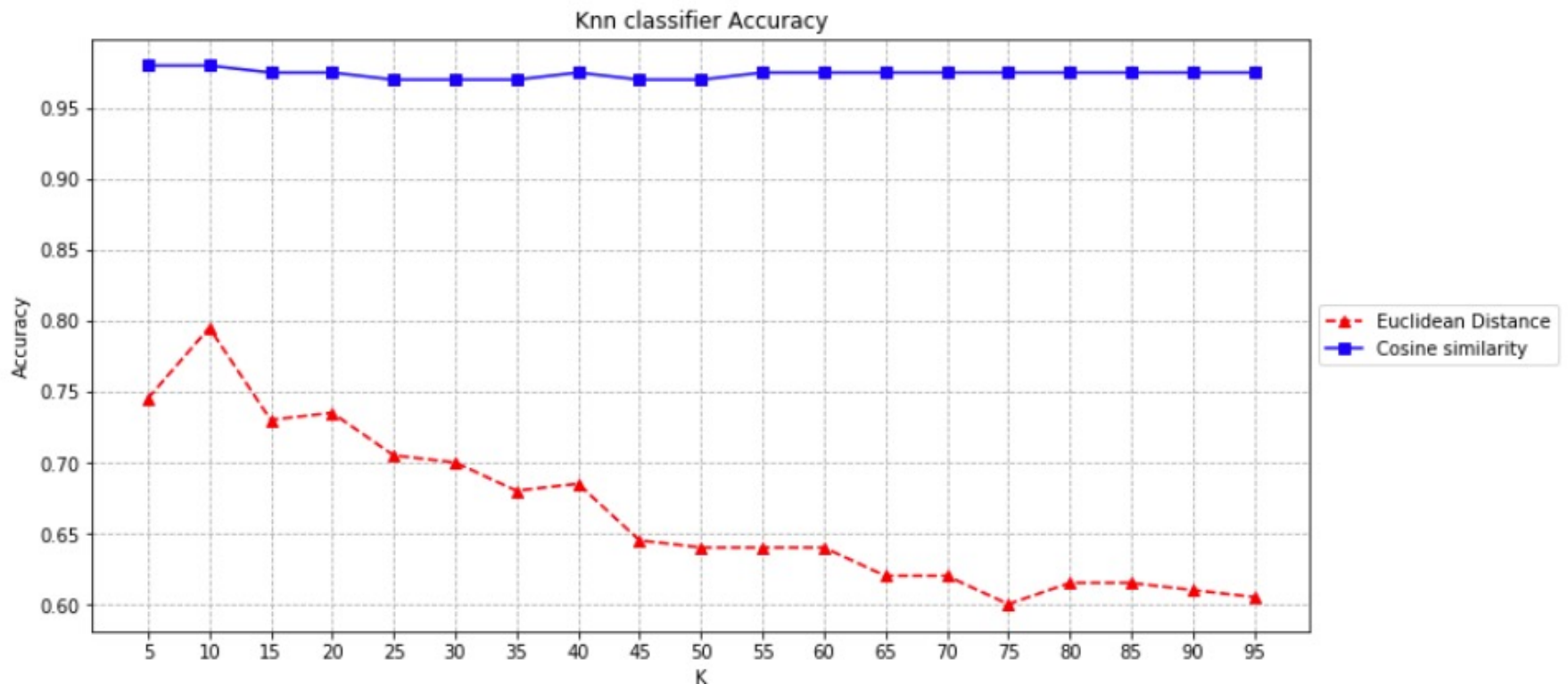
Wall time: 2min 36s

```
1 %%time  
2 for K in range(5, 100, 5):  
3     Cosine.append(knn_evaluate(DTM_test, test_lab_array, DTM_train, train_lab_array, K, 1))
```

Wall time: 52 s

Knn Classification & Evaluation

```
knn_evaluate(DT_test, test_labels, DT_train, train_labels, K, measure)
```



If you get stuck...

- Other questions depend on your KNN function
- If you get stuck, compare Euclidean and Cosine with scikit-learn's function

Rocchio Text Categorization Algorithm (Training)

Assume the set of categories is $\{c_1, c_2, \dots, c_n\}$

For i from 1 to n let $\mathbf{p}_i = \langle 0, 0, \dots, 0 \rangle$ (*init. prototype vectors*)

For each training example $\langle x, c(x) \rangle \in D$

Let \mathbf{d} be the TF/IDF term vector for doc x

Let $i = j$ where $c_j = c(x)$

(*sum all the document vectors in c_i to get \mathbf{p}_i*)

Let $\mathbf{p}_i = \mathbf{p}_i + \mathbf{d}$

Rocchio Text Categorization Algorithm (Test)

Given test document x

Let \mathbf{d} be the TF/IDF term vector for x

Let $m = -2$ (*init. maximum cosSim*)

For i from 1 to n :

(compute similarity to prototype vector)

Let $s = \text{cosSim}(\mathbf{d}, \mathbf{p}_i)$

if $s > m$

let $m = s$

let $r = c_i$ (*update most similar class prototype*)

Return class r

Rocchio Classification

```
def Rocchio_Train(train, labels):
```

```
    . . .
```

```
    return prototype
```

```
def Rocchio_classifier(prototype, instance):
```

```
    . . .
```

```
    return predicted_label, sims
```

```
def rocchio_evaluate(test, test_lab, prototype):
```

```
    . . .
```

```
    return accuracy
```

Note: **prototype** could be a dictionary with unique class labels as keys and one dimensional arrays representing the prototype vectors for the class as values

tf x idf Transformation

The initial
Term x Doc matrix
(Inverted Index)

	Doc 1	Doc 2	Doc 3	Doc 4	Doc 5	Doc 6
T1	0	2	4	0	1	0
T2	1	3	0	0	0	2
T3	0	1	0	2	0	0
T4	3	0	1	5	4	0
T5	0	4	0	0	0	1
T6	2	7	2	1	3	0
T7	1	0	0	5	5	1
T8	0	1	1	0	0	3

df	idf = $\log_2(N/df)$
3	1.00
3	1.00
2	1.58
4	0.58
2	1.58
5	0.26
4	0.58
3	1.00

Documents represented as vectors of words

tf x idf
Term x Doc matrix

	Doc 1	Doc 2	Doc 3	Doc 4	Doc 5	Doc 6
T1	0.00	2.00	4.00	0.00	1.00	0.00
T2	1.00	3.00	0.00	0.00	0.00	2.00
T3	0.00	1.58	0.00	3.17	0.00	0.00
T4	1.74	0.00	0.58	2.90	2.32	0.00
T5	0.00	6.34	0.00	0.00	0.00	1.58
T6	0.53	1.84	0.53	0.26	0.79	0.00
T7	0.58	0.00	0.00	2.92	2.92	0.58
T8	0.00	1.00	1.00	0.00	0.00	3.00