

Locality Sensitive Hashing

In this assignment you will implement locality sensitive hashing step-by-step and apply that on detecting almost duplicate names. You are required to implement the functions as instructed. Do not change the signatures of the functions.

Submission instruction

Submit a text file with extension .py containing the functions you have implemented. You can simply use the "download as" option from Jupyter notebook to export the notebook as a .py file. The filename should be LSH_{YOUR_ROLL_NO}.py where {YOUR_ROLL_NO} is your 6 letter roll number. For example, a sample file name would be LSH_CS1702.py. Please take this seriously and name the file correctly.

Important: do not change the function signatures

You need to implement the functions given in this notebook. You may define more functions at your will if you need to, but make sure you implement the functions given and do not change their signatures. Before submitting, remove unnecessary print / collect statements in order to avoid too much console output.

Setup environment

We need to setup pyspark and pydrive.

```
In [2]: 1 !pip install pyspark
        2 !pip install -U -q PyDrive
        3 !apt install openjdk-8-jdk-headless -qq
        4 import os
        5 os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
```

```
Requirement already satisfied: pyspark in /usr/local/lib/python3.6/dist-packages (3.0.1)
Requirement already satisfied: py4j==0.10.9 in /usr/local/lib/python3.6/dist-packages (from pyspark) (0.10.9)
openjdk-8-jdk-headless is already the newest version (8u272-b10-0ubuntu1~18.04).
0 upgraded, 0 newly installed, 0 to remove and 11 not upgraded.
```

Follow the interactive instructions for accessing the file in Google drive.

```
In [3]: 1 from pydrive.auth import GoogleAuth
        2 from pydrive.drive import GoogleDrive
        3 from google.colab import auth
        4 from oauth2client.client import GoogleCredentials
        5
        6 # Authenticate and create the PyDrive client
        7 auth.authenticate_user()
        8 gauth = GoogleAuth()
        9 gauth.credentials = GoogleCredentials.get_application_default()
       10 drive = GoogleDrive(gauth)
```

Load the appropriate file

The files are shared on Google drive, with the following ids given. Choose to keep the appropriate line, use the corresponding filename and comment the others. After executing, you should be able to see the file in the files section on the left panel.

You may face troubles trying to run your code on the full data, so it is totally okay to run it on one

of the smaller files.

```
In [4]: 1 #id='laoZykfz5GLGGw3lRA86ogd7yCsgkgHoT' # for titles-small.txt
2 id='1IPssUa3m-zfWmvVbLtI-lKhJbeZjIgBg' # for titles-10k.txt (with 10k titles)
3 # id='1RzTA4i0fH3b0iyG2kDaxweUuj4AEcMda' #
4 downloaded = drive.CreateFile({'id': id})
5 downloaded.GetContentFile('titles-10k.txt')
```

Parameters

Let us set our parameters. You are free to change these around and experiment. However, do not hardcode them anywhere else in the code.

```
In [5]: 1 #The following tunable parameters are considered for the assignment
2 # Number of hash functions
3 n = 20
4 # size of each shingle
5 k = 3
6 # Number of bands
7 h = 4
```

Note: Please open the notebook in Google colab otherwise it might show indentation error in local system because here 2 whitespace works for indentation but in local system, jupyter notebook takes 4 whitespaces or tab for indentation

1. Converting a string to shingles (of characters) [Marks: 5]

In our setup, the items we would compare are movie titles (strings). You would have to convert names to sets of k-shingles (of characters), after removing the whitespaces and converting to lowercase. For example, the 3-shingling of the name 'Die hard' would be ['die', 'ieh', 'eha', 'har', 'ard'].

Implement the following function which takes a string and k as input and outputs the list of unique shingles generated from the string.

```
In [6]: 1 def string_to_shingles(name, k):
2     shingles=[]
3     name=name.lower() # it converts given string to all lower characters string
4     string="".join(name.split()) # it removes the whitespaces from the string
5     #print(str)
6     for s in range(len(string)-k+1):
7         shingles.append(string[s:s+k]) # taking k-length shingles
8     return list(set(shingles)) # to get unique shingles, first convert list to set
9     # returns list of unique shingles generated from the string
```

Also test your code.

```
In [7]: 1 print(string_to_shingles('Die hard', 3))

['har', 'die', 'ard', 'ieh', 'eha']
```

2. Generating the shingle - item matrix from the whole data [Marks: 20]

For each title, first give it an ID. You can do it by the zipWithIndex() operation on Spark RDDs. Try testing it first.

Then, for each title, get the list of shingles in map and for each unique shingle, get the list of IDs using reduce.

Tips: when your input RDD is a list of tuples of the following form:

```
[('Die hard', 0), ('Die another day', 1), ('Tomorrow never dies', 2), ('Chup
ke chup ke', 3), .... ]
```

for each tuple `t`, you may use `t[0]`, `t[1]` etc to access the first, second (and so on) elements of the

The "titles-10k.txt" file has total 10k movie titles which are loaded in the titles RDD.

```
In [8]: 1 from pyspark import SparkContext, SparkConf
2 sc = SparkContext.getOrCreate()
3
4 # Use appropriate file path here
5 titles = sc.textFile("titles-10k.txt")
6
7 # Get the number of titles
8 N = titles.count() # N=10000 = total number of movie titles
9 itemsByShingles = titles.zipWithIndex() # associate a unique id for each ti
```

Here, we convert each title to a 3-shingles and keyval is a list of key-value pair in which key has title and value is a list of 3-shingles for corresponding title.

```
In [9]: 1 keyval=titles.map(lambda x: (x,string_to_shingles(x,k)))
2 #keyval.takeSample(False 5) # taking a sample of size 5 without replacement
```

Here, keyval1 returns the key-value pair where key is a list of 3-shingles for a particular title and corresponding movie id.

```
In [14]: 1 keyval1=itemsByShingles.map(lambda x: (string_to_shingles(x[0],k),x[1]))
2 #keyval1.takeSample(False 5)
```

From the list of " ([list of 3-shingles] , movie-id)", I make a pair of a particular shingle and movie id.

```
In [15]: 1 key_val=keyval1.flatMap(lambda n: [(x, n[1]) for x in n[0]])
2 #key_val.takeSample(False 5)
```

```
In [16]: 1 #print(key_val.count())
```

Here, result returns a key-value pair where key is a shingle and and value is the list of movie id associated with a particular shingle which is the objective of this part.

```
In [24]: 1 result=key_val.combineByKey(lambda v:[v],lambda x,y:x+[y],lambda x,y:x+y)
2 result.takeSample(False 5)
```

```
Out[24]: [('ulw', [988, 1451, 5609, 6125, 6898, 9356, 9504, 9602]),
('chä', [7889]),
('llw', [1054, 1951, 6024, 8702, 9575, 9732]),
('ayu', [2739, 3707]),
('m,t', [4006, 6058, 9482])]
```

```
In [25]: 1 #result.count() # returns total number of unique shingles
```

At the end of this step, the itemsByShingles RDD should have the list of (shingle, [list of movie ids]).

If you display your output, it should look like:

```
[(shingle, [list of movie ids]), (shingle, [list of movie ids]), ... ]
```

However, do not keep collect or large print statements in your code at the time of submitting.

3. Computation of min-hash signature matrix [Marks: 10 + 10 = 20]

Instead of using random permutations, you will implement min-hash function using Murmurhash (v3) functions, as discussed in the class. The input to your function should be a *title* (which corresponds to a row) and the output should be a number (hash value).

Recall the outline of the min-hash signature matrix algorithm:

1. For each row r BEGIN
2. Compute $h_1(r), h_2(r), \dots, h_n(r)$
3. For each column j BEGIN
4. If the j -th column has 1 in row r
5. For each $i = 1, 2, \dots, n$ BEGIN
6. set $SIG_{i,j} = \min(SIG_{i,j}, h_i(r))$
7. END
8. END IF
9. END
10. END

For each shingle, we have the list of movie IDs which contain the shingle. So, for each shingle (each row), we can perform the actions in lines 5-7 only for the movie IDs that are present in the list. For any other j , the i -th column does not have a 1 in row r .

Do not create a SIG matrix and try to update it. Instead, note that for every 1-entry in the shingle-id matrix (that is, every id in the list corresponding to a shingle), the corresponding (i, id) -th entry of the signature matrix may get potentially updated. Simply output the tuple $((i, id), h_i(\text{text}))$ in another map. The final value of the (i, id) -th entry of the signature matrix is the minimum of all such $h_i(\text{text})$ values obtained in this process. That minimum can be computed by another reduce process.

(a) The map [Marks: 10 out of 20]

To make it easier, you may implement an `update_signature` function which takes the shingle (text), the list of ids associated with the shingle and the total number of shingles (for computing the hash values). It should simply return the tuples $((i, id), h_i(\text{text}))$ for all $i = 1, \dots, n$.

Note: There has been a correction here. Instead of N (the number of movie ids), we need to pass the number of Shingles here.

In [27]: `1 !pip3 install mmh3`

Requirement already satisfied: mmh3 in /usr/local/lib/python3.6/dist-packages (2.5.1)

Here, mmh3 returns the key-value pair where key is a tuple of hash-function index and movie id and value is hash-value. Actually, here, key is $(i, j)^{th}$ entry in signature matrix and value is value in that $(i, j)^{th}$ entry for the signature matrix.

```
In [29]: 1 import mmh3
2 def update_signature(text, ids, numOfShingles):
3
4     # h = [mmh3.hash(text,i) % (total number of shingles) for i in range(n)]
5
6     min_sig = []
7     for i in range(n):
8         for j in range(len(ids)):
9             temp=[]
10            temp2=[]
11            temp.append(i)
12            temp.append(ids[j])
13            temp2.append(tuple(temp))
14            temp2.append(mmh3.hash(text,i) % (numOfShingles))
15            min_sig.append(tuple(temp2))
```

```

16
17
18 return min_sig

```

(b) The reduce [Marks: 10 out of 20]

Now that you have all $((i, id), h_i(\text{shingle}))$ tuples output by the map, use reduce to compute the minimum of $h_i(\text{shingle})$ for every (i, id) key. This would produce the signature matrix in a sparse matrix representation.

You may actually implement the map function above and call map and reduce together later as well.

Here, new_key_val returns the list of key-value pair where key is a pair of hash-function id $i = 0$ to $(n - 1)$ and a movie-id. Value in the key-value pair is value of the hash-function where input is a unique 3-shingle. Here, key-value pairs are corresponding to the unique 3-shingle and associated list of movie-ids.

```

In [31]: 1 numOfShingles = result.count() # Calculate, you need to pass this to the up
2
3 # You should use map with the update_signature function and a reduce
4 #new_key_val=result.flatMap(lambda n: [(x, n[1]) for x in n[1]])
5
6 new_key_val=result.map(lambda x: update_signature(x[0],x[1],numOfShingles)
7 #new_key_val.takeSample(False,2)
8 new_key_val = new_key_val.flatMap(lambda list: list)

```

```

In [32]: 1 #print(result.count(), new_key_val.count())

```

Now, signature returns the signature matrix which is our objective for this part. It is in the form of $((i, j), k)$ where, (i, j) represents the $(i, j)^{th}$ entry of the signature matrix and k is the value of the hash function.

Note: Here, size of the signature matrix is 199960, not $20 \times 10000 = 200000$. The reason is, in the original dataset "titles-10k.txt", there are 2 movie titles "99" and "S1" for which I can't generate 3-shingle. So, size of the signature matrix will be $20 * (10000 - 2) = 199960$.

```

In [33]: 1 signature = new_key_val.reduceByKey(min)
2 signature.takeSample(False, 5)

```

```

Out[33]: [((15, 7758), 1104),
          ((18, 6315), 899),
          ((4, 7382), 346),
          ((6, 4133), 525),
          ((17, 9220), 1166)]

```

```

In [34]: 1 #signature.count()

```

The RDD signature should be of the following form:

```

[((0, 1), 5), ((0, 3), 56), ... ]

```

where each tuple $((i, j), v)$ represents the (i, j) -th entry of the signature matrix with value v .

4. Implement the banding [Marks: 15]

Now configure your number of bands b (a divisor of number of hash functions n) and implement the candidate pair computation. Any pair of movies which agree completely in their signature on at least one band should be output as candidate pairs.

Here, signature matrix is divided into 4 parts by rows since there are 4 bands and number of rows in each band is 5.

```
In [35]: 1 #r=n/b = 20/4 = 5
2
3 band1=signature.filter(lambda x : x[0][0] >=0 and x[0][0]<=4)
4 band2=signature.filter(lambda x : x[0][0] >=5 and x[0][0]<=9)
5 band3=signature.filter(lambda x : x[0][0] >=10 and x[0][0]<=14)
6 band4=signature.filter(lambda x : x[0][0] >=15 and x[0][0]<=19)
7
8 candidate_pairs = []
```

Since, I get the candidate pairs in band1 where pair of movies completely agree in their signature, so I will not check the other bands. The logic to find the candidate pair which I have used is: first for a band1, I listed the repeating hash-values in the signature matrix along with their entry means to which row and column they belong to. After that, pair-wise columns in which hash-values are same, I compare the other 3 hash-values for those pair-wise columns. If all values are matches, then I return the those column pairs. **Note:** Here, I return only the candidate pairs of movie ids. Later, I will return candidate pairs by movie titles.

```
In [36]: 1 #band1
2
3 temp01= band1.map(lambda x : (x[1],x[0]))
4 temp101=temp01.groupByKey().map(lambda x : (x[0], list(x[1])))
5 #print(temp101.takeSample(False,2))
6 def ls(x):
7     f2=[]
8     for i in range(len(x[1])):
9         for j in range(i+1,len(x[1])):
10             f1=[]
11             if x[1][i][0]== x[1][j][0]:
12                 l1=[]
13                 l1.append(x[1][i][1])
14                 l1.append(x[1][j][1])
15                 f1.append(tuple(l1))
16                 f1.append(x[1][j][0])
17                 #f1.append(x[0])
18                 f2.append(tuple(f1))
19     return f2
20 filt1=temp101.flatMap(lambda x : [ls(x)])
21 #filt1.take(2)
22 filt1=filt1.flatMap(lambda l: l)
23 #filt1.count()
24 fin1=filt1.groupByKey().map(lambda x : (x[0], list(x[1])))
25 #print(fin1.takeSample(False,2))
26 #fin1.count()
27 def check(x):
28     if len(x[1])==4:
29         return x[0]
30 res1=fin1.map(lambda x: check(x))
31 res1=res1.filter(lambda x: x is not None)
32 res1.collect()
```

Out[36]:

```
[(7456, 9558),
 (3374, 988),
 (3374, 4620)]
```

Now, instead of candidate pairs as movie ids, I will return the candidate pairs by movie titles. I used the above code (just copy-pasted) and instead of movie-id, I replaced the movie-titles.

```
In [37]: 1 keyval_modified=keyval.map(lambda x: (string_to_shingles(x[0],k),x[0])) #it
2 #list of 3-shingles for a particular title and corresponding movie id
3 #keyval_modified.takeSample(False,5)
4 key_val_modified=keyval_modified.flatMap(lambda n: [(x, n[1]) for x in n[0]]
5 # a pair of a particular shingle and movie id.
6 #key_val_modified.takeSample(False,5)
7 result_modified=key_val_modified.combineByKey(lambda v:[v],lambda x,y:x+[y]
8 # where key is a shingle and value is the list of movie id associated w
9 #result_modified.takeSample(False,5)
10 #numOfShingles = result.count() # Calculate, you need to pass this to the u
11
12 #signature = # Implement the rest here
13
14 # You should use map with the update_signature function and a reduce
15 #new_key_val=result.flatMap(lambda n: [(x, n[1]) for x in n[1]])
16
17 new_key_val_modified=result_modified.map(lambda x: update_signature(x[0],x[
18 #new_key_val_modified.takeSample(False,2)
19 #signature.take(5)
20 new_key_val_modified = new_key_val_modified.flatMap(lambda list: list)
21 signature_modified = new_key_val_modified.reduceByKey(min)
22 #r=n/b = 20/4 = 5
23 #x = sc.parallelize([1,2,4,5])
24 ##max_val = x.reduce(lambda a, b: a if a > b else b)
25 #print(max_val)
26 band1_modified=signature_modified.filter(lambda x : x[0][0] >=0 and x[0][0]
27 band2_modified=signature_modified.filter(lambda x : x[0][0] >=5 and x[0][0]
28 band3_modified=signature_modified.filter(lambda x : x[0][0] >=10 and x[0][0]
29 band4_modified=signature_modified.filter(lambda x : x[0][0] >=15 and x[0][0]
30
31 candidate_pairs = []
32 #band1_modified
33
34 temp01_modified= band1_modified.map(lambda x : (x[1],x[0]))
35 temp101_modified=temp01_modified.groupByKey().map(lambda x : (x[0], list(x[
36 #print(temp101.takeSample(False,2))
37 def ls(x):
38     f2=[]
39     for i in range(len(x[1])):
40         for j in range(i+1,len(x[1])):
41             f1=[]
42             if x[1][i][0]== x[1][j][0]:
43                 l1=[]
44                 l1.append(x[1][i][1])
45                 l1.append(x[1][j][1])
46                 f1.append(tuple(l1))
47                 f1.append(x[1][j][0])
48                 #f1.append(x[0])
49                 f2.append(tuple(f1))
50     return f2
51 filt1_modified=temp101_modified.flatMap(lambda x : [ls(x)])
52 #filt1.take(2)
53 filt1_modified=filt1_modified.flatMap(lambda l: l)
54 #filt1.count()
55 fin1_modified=filt1_modified.groupByKey().map(lambda x : (x[0], list(x[1])))
56 #print(fin1.takeSample(False,2))
57 #fin1.count()
58 def check(x):
59     if len(x[1])==4:
60         return x[0]
61 res1_modified=fin1_modified.map(lambda x: check(x))
```

```
Out[37]: [('Statistinnen des Lebens', 'Paul Wang's Destiny'),
          ('Statistinnen des Lebens', 'The Thread of Destiny'),
          ('Paul Wang's Destiny', 'The Thread of Destiny'),
          ('Caught in the Rain', 'Caught in a Park'),
          ('The Fighting Brothers', 'The Fighting Hope'),
          ('The Fighting Brothers', 'Fighting Mad'),
          ('The Fighting Brothers', 'The Fighting Trail'),
          ('The Fighting Hope', 'Fighting Mad'),
          ('Fighting Mad', 'The Fighting Trail'),
          ('Threads of Destiny', 'Destiny'),
          ('What's Sauce for the Goose', 'That's Him'),
          ('His Ups and Downs', 'Ups and Downs'),
          ('Hands Down', 'Pines Up and Palms Down'),
          ('The Widow's Secret', 'The Widow's Kids'),
          ('In the Window Recess', 'The Face at the Window'),
          ('The Governor's Ghost', 'The Governor's Lady'),
          ('The Joneses Have Amateur Theatricals', 'In the Theatrical Business'),
          ('His Prehistoric Past', 'Prehistoric Poultry'),
          ('The Price He Paid', 'The Price She Paid'),
          ('Paul's Prehistoric African Hunt', 'Prehistoric African Hunt')]
```

5. Optional:

In []: 1