

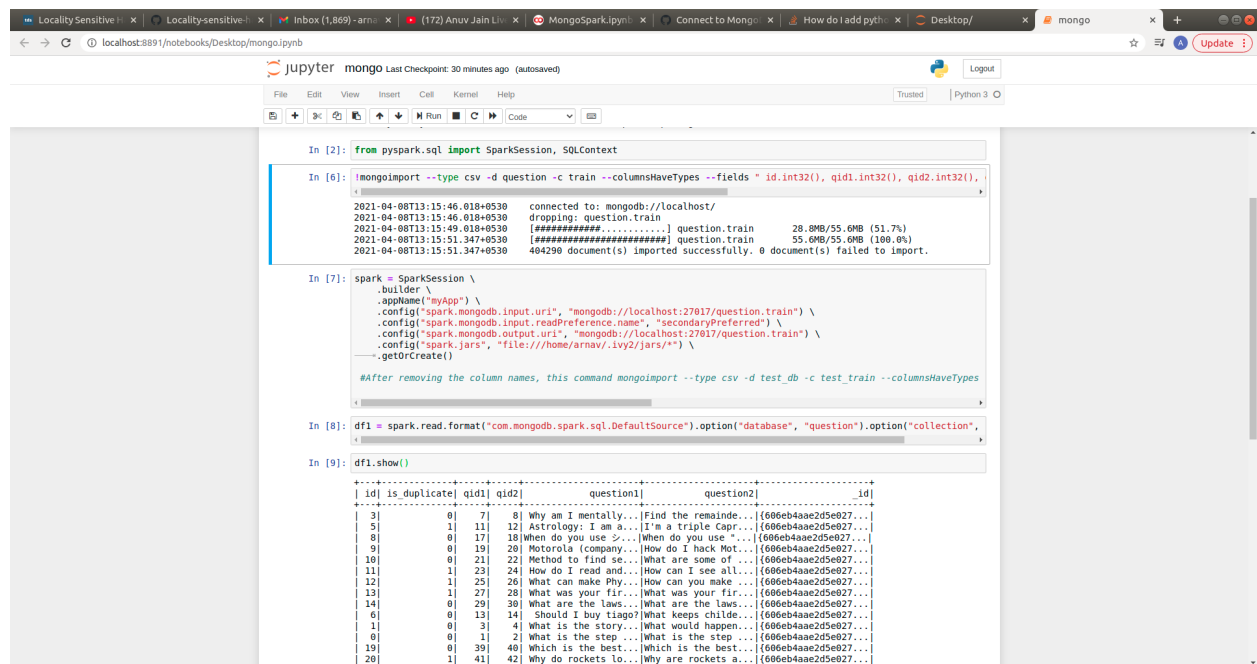
LSH using PySpark

Garvita Jain 2018034 | Arnav Tandon 2018278

Store data in MongoDB

Command : `!mongoimport --type csv -d question -c train --columnsHaveTypes --fields "id.int32(), qid1.int32(), qid2.int32(), question1.string(), question2.string(), is_duplicate.int32()" --drop train.csv`

Output



```
In [2]: from pyspark.sql import SparkSession, SQLContext

In [6]: !mongoimport --type csv -d question -c train --columnsHaveTypes --fields "id.int32(), qid1.int32(), qid2.int32(),
2021-04-08T13:15:46.018+0530 connected to: mongodb://localhost/
2021-04-08T13:15:46.018+0530 dropping: question.train
2021-04-08T13:15:49.018+0530 [#####] question.train 28.6MB/55.6MB (51.7%)
2021-04-08T13:15:51.347+0530 [#####] question.train 55.6MB/55.6MB (100.0%)
2021-04-08T13:15:51.347+0530 404290 document(s) imported successfully. 0 document(s) failed to import.

In [7]: spark = SparkSession \
        .builder \
        .appName("myApp") \
        .config("spark.mongodb.input.uri", "mongodb://localhost:27017/question.train") \
        .config("spark.mongodb.input.readPreference.name", "secondaryPreferred") \
        .config("spark.mongodb.output.uri", "mongodb://localhost:27017/question.train") \
        .config("spark.jars", "file:///home/arnav/.ivy2/jars/*") \
        .getOrCreate()

#After removing the column names, this command mongoimport --type csv -d test_db -c test_train --columnsHaveTypes

In [8]: df1 = spark.read.format("com.mongodb.spark.sql.DefaultSource").option("database", "question").option("collection",
df1.show()
```

id	is_duplicate	qid1	qid2	question1	question2	_id
3	0	7	8	Why am I mentally...Find the remaine...	I'm a triple Capr...	{606eb4aae2d5e027...}
5	1	11	12	Astrology: I am a...		{606eb4aae2d5e027...}
8	0	17	18	When do you use S...	When do you use "...	{606eb4aae2d5e027...}
9	0	19	20	Motorola (company)...	How do I hack Mot...	{606eb4aae2d5e027...}
10	0	21	22	Method to find se...	What are some of ...	{606eb4aae2d5e027...}
11	1	23	24	How do I read and...	How can I see all...	{606eb4aae2d5e027...}
12	1	25	26	What can make Phy...	How can you make ...	{606eb4aae2d5e027...}
13	1	27	28	What was your fir...	What was your fir...	{606eb4aae2d5e027...}
14	0	29	30	What are the laws...	What are the laws...	{606eb4aae2d5e027...}
6	0	13	14	Should I buy tiago?	What keeps childe...	{606eb4aae2d5e027...}
1	0	3	4	What is the story...	What would happen...	{606eb4aae2d5e027...}
0	0	1	2	What is the step ...	What is the step ...	{606eb4aae2d5e027...}
19	0	39	40	Which is the best...	Which is the best...	{606eb4aae2d5e027...}
20	1	41	42	Why do rockets lo...	Why are rockets a...	{606eb4aae2d5e027...}

Reading data normally :

```
df1 = spark.read.format("com.mongodb.spark.sql.DefaultSource").option("database",
"question").option("collection", "train").load()
```

```
df2_rdd = spark.read.format("com.mongodb.spark.sql.DefaultSource").option("database",
"question").option("collection", "train").load().rdd
```

jupyter mongo Last Checkpoint: 35 minutes ago (unsaved changes)

File Edit View Insert Cell Kernel Help Trusted Python 3

```
In [8]: tSource").option("database", "question").option("collection", "train").load().option("collection", "train").load()

In [9]: df1.show()
```

	id	is_duplicate	qid1	qid2	question1	question2	_id
3		0	7	8	Why am I mentally...	Find the remainde...	{606eb4aae2d5e027...
5		1	11	12	Astrology: I am a...	I'm a triple Capr...	{606eb4aae2d5e027...
8		0	17	18	When do you use シ...	When do you use "...	{606eb4aae2d5e027...
9		0	19	20	Motorola (company...	How do I hack Mot...	{606eb4aae2d5e027...
10		0	21	22	Method to find se...	What are some of ...	{606eb4aae2d5e027...
11		1	23	24	How do I read and...	How can I see all...	{606eb4aae2d5e027...
12		1	25	26	What can make Phy...	How can you make ...	{606eb4aae2d5e027...
13		1	27	28	What was your fir...	What was your fir...	{606eb4aae2d5e027...
14		0	29	30	What are the laws...	What are the laws...	{606eb4aae2d5e027...
6		0	13	14	Should I buy tiago?	What keeps child...	{606eb4aae2d5e027...
1		0	3	4	What is the story...	What would happen...	{606eb4aae2d5e027...
0		0	1	2	What is the step ...	What is the step ...	{606eb4aae2d5e027...
19		0	39	40	Which is the best...	Which is the best...	{606eb4aae2d5e027...
20		1	41	42	Why do rockets lo...	Why are rockets a...	{606eb4aae2d5e027...
21		0	43	44	What's causing so...	What can I do to ...	{606eb4aae2d5e027...
22		0	45	46	What are the ques...	Which question sh...	{606eb4aae2d5e027...
23		0	47	48	How much is 30 kV...	Where can I find ...	{606eb4aae2d5e027...
24		0	49	50	What does it mean...	How many times a ...	{606eb4aae2d5e027...
25		0	51	52	What are some tip...	What are some tip...	{606eb4aae2d5e027...
26		0	53	54	What is web appli...	What is the web a...	{606eb4aae2d5e027...

only showing top 20 rows

```
In [12]: df2_rdd = spark.read.format("com.mongodb.spark.sql.DefaultSource").option("database", "question").option("collection", "train").load()

In [13]: df2_rdd.collect()
```

```
Out[13]: [Row(id=3, is_duplicate=0, qid1=7, qid2=8, question1='Why am I mentally very lonely? How can I solve it?', question2='Find the remainder when [math]23^{24}/[math] is divided by 24,23', _id=Row(oid='606eb4aae2d5e027f20d4d31')),
Row(id=5, is_duplicate=1, qid1=11, qid2=12, question1='Astrology: I am a Capricorn Sun Cap moon and cap ri sing...what does that say about me?', question2='I'm a triple Capricorn (Sun, Moon and ascendant in Capricorn) What does this say about me?', _id=Row(oid='606eb4aae2d5e027f20d4d32')),
Row(id=8, is_duplicate=0, qid1=17, qid2=18, question1='When do you use シ instead of L?', question2='When do you use "&" instead of "and"?', _id=Row(oid='606eb4aae2d5e027f20d4d33')),
Row(id=9, is_duplicate=0, qid1=19, qid2=20, question1='Motorola (company): Can I hack my Charter Motorola DCX3400?', question2='How do I hack Motorola DCX3400 for free internet?', _id=Row(oid='606eb4aae2d5e027f20d4d34'))]
```

localhost:8889/notebooks/Downloads/code.ipynb

jupyter code (autosaved)

File Edit View Insert Cell Kernel Help Connecting to kernel Trusted Python 3

```
In [3]: spark = SparkSession \
    .builder \
    .appName("Assignment3") \
    .config("spark.mongodb.input.uri", "mongodb://localhost:27017/question.train") \
    .config("spark.mongodb.output.uri", "mongodb://localhost:27017/question.train") \
    .config("spark.jars", "file:///home/arnav/.ivy2/jars/**") \
    .getOrCreate()

In [4]: !mongoimport --type csv -d question -c train --headerline --maintainInsertionOrder --drop train.csv
```

```
2021-04-09T20:16:13.067+0530 connected to: mongodb://localhost/
2021-04-09T20:16:13.067+0530 dropping: question.train
2021-04-09T20:16:16.067+0530 [#####] question.train 33.1MB/60.5MB (54.7%)
2021-04-09T20:16:19.067+0530 [#####] question.train 59.0MB/60.5MB (97.5%)
2021-04-09T20:16:19.213+0530 [#####] question.train 60.5MB/60.5MB (100.0%)
2021-04-09T20:16:19.213+0530 404290 document(s) imported successfully. 0 document(s) failed to import.

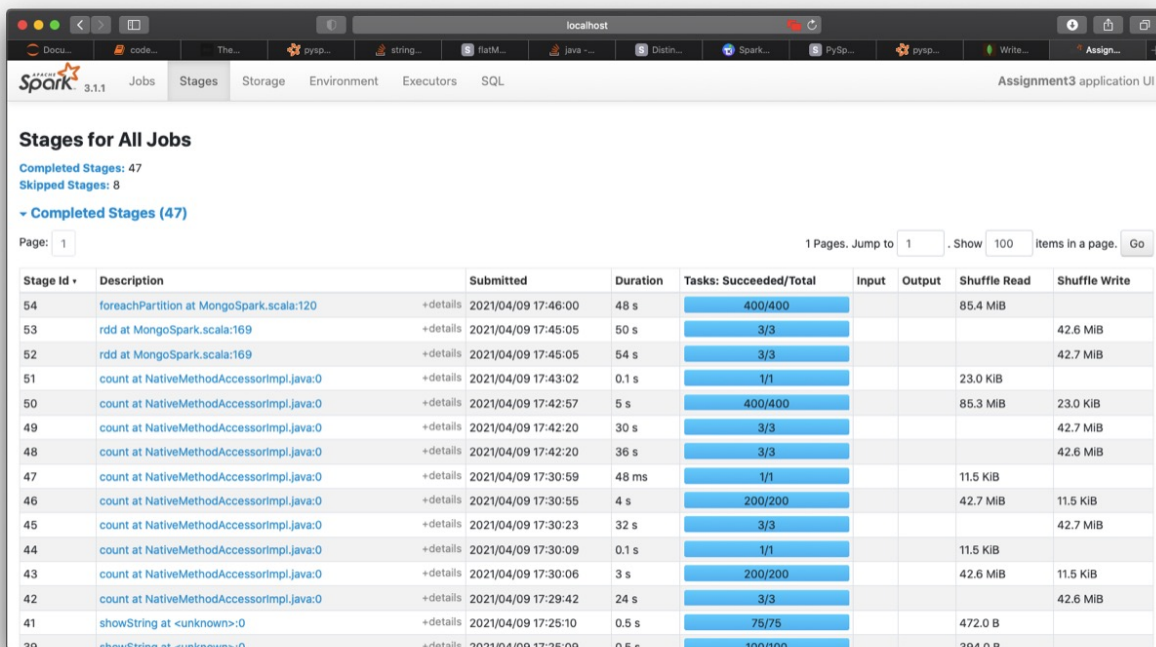
In [5]: spark.sparkContext.uiWebUrl
Out[5]: 'http://192.168.1.43:4041'

In [6]: df1 = spark.read.format("com.mongodb.spark.sql.DefaultSource").option("database", "question") \
    .option("collection", "train").load().limit(100000)

In [7]: # df1.count()
df1.show()
```

	id	is_duplicate	qid1	qid2	question1	question2
{607068b589c0716f...}	0	0	1	2	What is the step ...	What is the step ...
{607068b589c0716f...}	1	0	3	4	What is the story...	What would happen...
{607068b589c0716f...}	2	0	5	6	How can I increas...	How can Internet ...
{607068b589c0716f...}	3	0	7	8	Why am I mentally...	Find the remainde...
{607068b589c0716f...}	4	0	9	10	Which one dissolv...	Which fish would ...
{607068b589c0716f...}	5	1	11	12	Astrology: I am a...	I'm a triple Capr...
{607068b589c0716f...}	6	0	13	14	Should I buy tiago?	What keeps child...
{607068b589c0716f...}	7	1	15	16	How can I be a go...	What should I do ...
{607068b589c0716f...}	8	0	17	18	When do you use シ...	When do you use "...
{607068b589c0716f...}	9	0	19	20	Motorola (company...	How do I hack Mot...
{607068b589c0716f...}	10	0	21	22	Method to find se...	What are some of ...
{607068b589c0716f...}	11	1	23	24	How do I read and...	How can I see all...
{607068b589c0716f...}	12	1	25	26	What can make Phy...	How can you make ...
{607068b589c0716f...}	13	1	27	28	What was your fir...	What was your fir...
{607068b589c0716f...}	14	0	29	30	What are the laws...	What are the laws...
{607068b589c0716f...}	15	1	31	32	What would a Trum...	How will a Trump ...

Apache Spark Web UI screenshots :



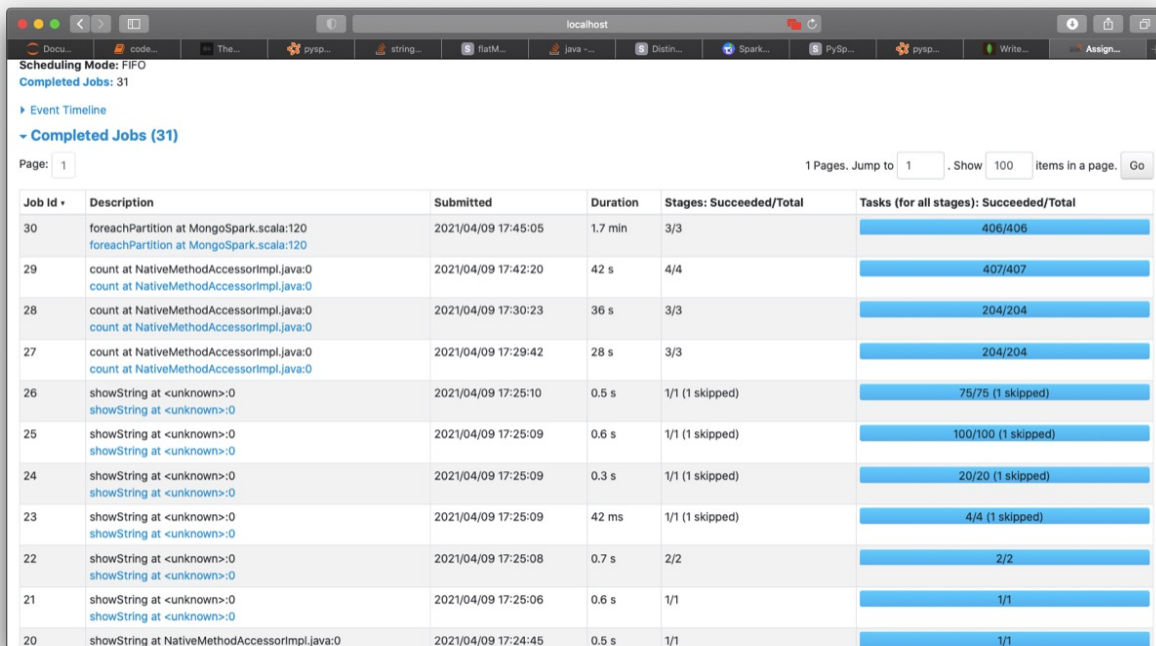
Stages for All Jobs

Completed Stages: 47
Skipped Stages: 8

▼ Completed Stages (47)

Page: 1 1 Pages. Jump to 1 . Show 100 Items in a page. Go

Stage ID	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
54	foreachPartition at MongoSpark.scala:120	2021/04/09 17:46:00	48 s	400/400			85.4 MiB	
53	rdd at MongoSpark.scala:169	2021/04/09 17:45:05	50 s	3/3				42.6 MiB
52	rdd at MongoSpark.scala:169	2021/04/09 17:45:05	54 s	3/3				42.7 MiB
51	count at NativeMethodAccessorImpl.java:0	2021/04/09 17:43:02	0.1 s	1/1			23.0 KiB	
50	count at NativeMethodAccessorImpl.java:0	2021/04/09 17:42:57	5 s	400/400			85.3 MiB	23.0 KiB
49	count at NativeMethodAccessorImpl.java:0	2021/04/09 17:42:20	30 s	3/3				42.7 MiB
48	count at NativeMethodAccessorImpl.java:0	2021/04/09 17:42:20	36 s	3/3				42.6 MiB
47	count at NativeMethodAccessorImpl.java:0	2021/04/09 17:30:59	48 ms	1/1			11.5 KiB	
46	count at NativeMethodAccessorImpl.java:0	2021/04/09 17:30:55	4 s	200/200			42.7 MiB	11.5 KiB
45	count at NativeMethodAccessorImpl.java:0	2021/04/09 17:30:23	32 s	3/3				42.7 MiB
44	count at NativeMethodAccessorImpl.java:0	2021/04/09 17:30:09	0.1 s	1/1			11.5 KiB	
43	count at NativeMethodAccessorImpl.java:0	2021/04/09 17:30:06	3 s	200/200			42.6 MiB	11.5 KiB
42	count at NativeMethodAccessorImpl.java:0	2021/04/09 17:29:42	24 s	3/3				42.6 MiB
41	showString at <unknown>:0	2021/04/09 17:25:10	0.5 s	75/75			472.0 B	
39	showString at <unknown>:0	2021/04/09 17:25:09	0.5 s	100/100			394.0 B	



Scheduling Mode: FIFO

Completed Jobs: 31

▼ Completed Jobs (31)

Page: 1 1 Pages. Jump to 1 . Show 100 Items in a page. Go

Job ID	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
30	foreachPartition at MongoSpark.scala:120 foreachPartition at MongoSpark.scala:120	2021/04/09 17:45:05	1.7 min	3/3	406/406
29	count at NativeMethodAccessorImpl.java:0 count at NativeMethodAccessorImpl.java:0	2021/04/09 17:42:20	42 s	4/4	407/407
28	count at NativeMethodAccessorImpl.java:0 count at NativeMethodAccessorImpl.java:0	2021/04/09 17:30:23	36 s	3/3	204/204
27	count at NativeMethodAccessorImpl.java:0 count at NativeMethodAccessorImpl.java:0	2021/04/09 17:29:42	28 s	3/3	204/204
26	showString at <unknown>:0 showString at <unknown>:0	2021/04/09 17:25:10	0.5 s	1/1 (1 skipped)	75/75 (1 skipped)
25	showString at <unknown>:0 showString at <unknown>:0	2021/04/09 17:25:09	0.6 s	1/1 (1 skipped)	100/100 (1 skipped)
24	showString at <unknown>:0 showString at <unknown>:0	2021/04/09 17:25:09	0.3 s	1/1 (1 skipped)	20/20 (1 skipped)
23	showString at <unknown>:0 showString at <unknown>:0	2021/04/09 17:25:09	42 ms	1/1 (1 skipped)	4/4 (1 skipped)
22	showString at <unknown>:0 showString at <unknown>:0	2021/04/09 17:25:08	0.7 s	2/2	2/2
21	showString at <unknown>:0 showString at <unknown>:0	2021/04/09 17:25:06	0.6 s	1/1	1/1
20	showString at NativeMethodAccessorImpl.java:0	2021/04/09 17:24:45	0.5 s	1/1	1/1

Preprocessing the data :

- 1) We have cleaned the data by first lowering all the document content and then by using `regexp_replace` for removal of integers and symbols.
- 2) We have then tokenized the data on the dataframe which was obtained after cleaning the data in the above step.

- 3) Stopwords are removed from the pyspark.ml.feature library to remove the unwanted words which do not contribute significantly to the text obtained.

Shingling :

- 1) Shingles were created by splitting by whitespaces.
- 2) For each document, a unique set of shingles was created and added to a new table.

The screenshot shows a Jupyter Notebook with a code cell containing the following PySpark code:

```
In [12]: cv = CountVectorizer(inputCol="q_clean", outputCol="features", minDF=2.0)
model = cv.fit(df_final)
result = model.transform(df_final)
result.show(truncate=False)
```

The output of the code is a table with two columns: 'id' and 'features'. The 'id' column contains document IDs, and the 'features' column contains a list of shingles (words or phrases) for each document. The shingles are represented as a sparse matrix of 1s and 0s, where 1 indicates the presence of a shingle and 0 indicates its absence.

id	features
0	[step, step, guide, invest, share, market, india]
449,590,1223,2454	[1.0,1.0,1.0,1.0,2.0,1.0]
1	[story, kohinoor, kohinoor, diamond]
43,9550	[1.0,1.0,2.0]
2	[increase, speed, internet, connection, using, vpn]
7,247,323,1889,2080	[1.0,1.0,1.0,1.0,1.0]
3	[mentally, lonely, solve]
20,2532	[1.0,1.0,1.0]
14	[one, dissolve, water, quickly, sugar, salt, methane, carbon, di, oxide]
1586,1756,1992,6763,7113,9115,10030	[1.0,1.0,1.0,1.0,1.0,1.0,1.0]
15	[astrology, capricorn, sun, cap, moon, cap, risingwhat, say]
6,771,2270,4609,7650	[1.0,1.0,1.0,1.0,2.0,1.0]
6	[buy, tiago]
591	[1.0,1.0]
17	[good, geologist]
11	[1.0,1.0]
8	[use, , instead]
60	[1.0,1.0,1.0]
9	[motorola, company, hack, charter, motorolla, dcx]
3,6892,6994,21828	[1.0,1.0,1.0,1.0,1.0]
10	[method, find, separation, slits, using, fresnel, biprism]
846,3780	[1.0,1.0,1.0,1.0]
11	[read, find, youtube, comments]
2,172,1559	[1.0,1.0,1.0,1.0]
12	[make, physics, easy, learn]
35,536	[1.0,1.0,1.0,1.0]
113	[first, sexual, experience, like]
81,1321	[1.0,1.0,1.0,1.0]
14	[laws, change, status, student, visa, green, card, us, compare, immigration, laws, canada]
99,126,139,380,415,508,732,1166,1659	[1.0,1.0,1.0,1.0,1.0,1.0,1.0,2.0,1.0,1.0]
15	[trump, presidency, mean, current, international, masters, students, f, visa]
184,377,415,442,636,1030,1523	[1.0,1.0,1.0,1.0,1.0,1.0,1.0,1.0,1.0]
16	[manipulation, mean]
12	[1.0,1.0]

Sparse Vector of the document :

- 1) For each document, a sparse matrix was created to represent the presence/absence of each shingle present in the entire dataset.
- 2) The SparseVector provided the number of shingles/rows in the vector, the list of indices at which the shingle was present for a particular question, and its frequency.

```

12 | 15 | [increase, speed, internet, connection, using, vpn] | (27526, [71, 11]
7,247,322,1888,2088], [1.0,1.0,1.0,1.0,1.0,1.0]) |
13 | 7 | [mentally, lonely, solve] | (27526, [533, 25
20,2532], [1.0,1.0,1.0]) |
14 | 19 | [one, dissolve, water, quickly, sugar, salt, methane, carbon, di, oxide] | (27526, [7, 136,
1586,1756,1992,6763,7113,9115,10030], [1.0,1.0,1.0,1.0,1.0,1.0,1.0,1.0,1.0,1.0]) |
15 | 11 | [astrology, capricorn, sun, cap, moon, cap, risingwhat, say] | (27526, [137, 75
6,771,2270,4005,7050], [1.0,1.0,1.0,1.0,2.0,1.0]) |
16 | 13 | [buy, tiago] | (27526, [55, 187
59], [1.0,1.0]) |
17 | 15 | [good, geologist] | (27526, [6, 1826
11], [1.0,1.0]) |
18 | 17 | [use, , instead] | (27526, [0, 18, 3
60], [1.0,1.0,1.0]) |
19 | 19 | [motorola, company, hack, charter, motorola, dcx] | (27526, [98, 31
3,6892,6994,21828], [1.0,1.0,1.0,1.0,1.0]) |
10 | 21 | [method, find, separation, slits, using, fresnel, biprism] | (27526, [29, 71,
846, 3798], [1.0,1.0,1.0,1.0]) |
11 | 23 | [read, find, youtube, comments] | (27526, [29, 12
2,172,1559], [1.0,1.0,1.0,1.0]) |
12 | 25 | [make, physics, easy, learn] | (27526, [9, 15, 5
35,536], [1.0,1.0,1.0,1.0]) |
13 | 27 | [first, sexual, experience, like] | (27526, [5, 35, 1
81,1321], [1.0,1.0,1.0,1.0]) |
14 | 29 | [laws, change, status, student, visa, green, card, us, compare, immigration, laws, canada] | (27526, [24, 86,
99,126,139,380,415,508,732,1166,1659], [1.0,1.0,1.0,1.0,1.0,1.0,1.0,1.0,2.0,1.0,1.0]) |
15 | 31 | [trump, presidency, mean, current, international, masters, students, f, visa] | (27526, [34, 38,
184,377,415,442,636,1030,1523], [1.0,1.0,1.0,1.0,1.0,1.0,1.0,1.0,1.0,1.0]) |
16 | 33 | [manipulation, mean] | (27526, [38, 177
12], [1.0,1.0]) |
17 | 35 | [girls, want, friends, guy, reject] | (27526, [37, 15
3,160,165,3708], [1.0,1.0,1.0,1.0,1.0]) |
18 | 37 | [many, quora, users, posting, questions, readily, answered, google] | (27526, [10, 21,
47,63,662,1076,3266,9901], [1.0,1.0,1.0,1.0,1.0,1.0,1.0,1.0,1.0,1.0]) |
19 | 39 | [best, digital, marketing, institution, banglore] | (27526, [1, 246,
347,4492,7490], [1.0,1.0,1.0,1.0,1.0]) |
-----
only showing top 20 rows

```

```

In [13]: # result.first()

In [14]: '''
a hash function can be randomly assigned as ax+b*no of shingles where x=shingle_index of shingle
present in a document dx. This will give the value of hash function for that particular shingle_index.
We have to find the value corresponding to all shingles present in the document.
'''

```

Made Signature Matrix using random Hash functions :

We have taken $h = 30$ which are our number of hash functions and we have defined our hash function as $val = (a*i+b) \% (\text{length of the shingles})$ where a, b are random numbers.

A hash function can be randomly assigned as $ax+b\%no_of_shingles$ where $x=shingle_index$ of shingle present in a document dx . This will give the value of hash function for that particular shingle_index. We have to find the value corresponding to all shingles present in the document.

```

In [14]: '''
a hash function can be randomly assigned as ax+b*no of shingles where x=shingle_index of shingle
present in a document dx. This will give the value of hash function for that particular shingle_index.
We have to find the value corresponding to all shingles present in the document.
'''

Out[14]: '\na hash function can be randomly assigned as ax+b*no of shingles where x=shingle_index of shingle\n
present in a document dx. This will give the value of hash function for that particular shingle_index. \n
We have to find the value corresponding to all shingles present in the document. \n'

In [15]: h = 40 #no of hash functions
values_a = rnd.sample(range(1,2000), h)
values_b = rnd.sample(range(1,2000), h)

def minhash(a, b, ishingles_in_a_doc, len_shingles):
    hashes = []
    for i in ishingles_in_a_doc:
        val = (a+i*b)%len_shingles
        hashes.append(val)
    if hashes:
        return min(hashes)
    else:
        return 1000000

def create_sign_doc(document):
    docid = document['id'], document['qid']
    shingleids = document['features'].indices
    signature = []
    for a,b in zip(values_a, values_b):
        signature.append(minhash(a,b,shingleids, 54153))
    return((docid, signature))

In [16]: create_sign_doc(Row(id=0, qid=1, q_clean=['step', 'step', 'guide', 'invest', 'share', 'market', 'india'],
features=SparseVector(54153, {3: 1.0, 265: 1.0, 445: 1.0, 627: 1.0, 1153: 2.0, 2382: 1.0})))

Out[16]: ((0, 1),
[3501,
638,
5251,
4127,
1405,
4564,
3931,
1943,
374,
2164,
2068,
6868])

```

LSH Implementation

When we have to compare two things and the size is too it takes $O(n^2)$ time, but there is a better way of comparing every possible pair by **Locality Sensitive Hashing** (LSH) in $O(n)$.

As there are 2 lac short documents for comparison, so we are taking k as 5 (commonly used size for shorter documents).

We use the band structure on the signature matrix made by us earlier, An if there are n rows, so it into b bands and with each row having width w, so we get :

- 1) $N = b*w$
- 2) Let's say the true similarity score is p when comparison is done between a pair,
Probability of one band not matching = $1-p^w$

Probability of number of bands that match = $(1-p^w)^b$

So Probability of at least one band matches = $1-(1-p^w)^b$

After implementing the LSH in PySpark we got the candidate pairs, meaning they have a similarity score above the threshold and there are multiple candidate pairs with same band_hash as shown in the figure below.

B and W make sense here and they do the opposite things as when we increase b it gives documents more chances to match, so we get pairs with lower similarity scores in the candidate pair list. But increasing W makes the match criteria stricter, so it is restricted to higher similarity scores.

```
In [89]: 1 lsh_df.groupby('band_no', 'band_hash').agg(collect_list("doc_id").alias("candidate_pairs")).show()
```

band_no	band_hash	candidate_pairs
0	-9098964674786157832	[{1013, 2021}]
0	-9097076610132250771	[{92908, 155452}, ...]
0	-9075903486970713091	[{11085, 21433}, ...]
0	-9075133033449877333	[{45196, 81020}]
0	-9063007082634758990	[{94866, 158331}]
0	-9036711468021138188	[{463, 924}]
0	-8991721065649749160	[{59100, 103569}]
0	-8950445860725111188	[{98166, 163177}]
0	-8847112719701112263	[{97503, 162224}]
0	-8840263500075329849	[{76388, 130589}]
0	-8837973820842852960	[{28236, 52378}]
0	-8820194762220627155	[{33423, 61408}, ...]
0	-8767193704393368565	[{83793, 141777}]
0	-8733549606059257863	[{11752, 22682}, ...]
0	-8724898006211925323	[{87139, 146804}]
0	-8723675710942797727	[{12071, 23284}]
0	-8709671219524383574	[{84195, 142372}]
0	-8701223129789196813	[{29053, 53796}]
0	-8686888307652572036	[{11143, 21544}]
0	-8671445867209440873	[{73614, 126310}]

only showing top 20 rows

Precision, recall and confusion matrix :

1. TN / True Negative: case was negative and predicted negative
2. TP / True Positive: case was positive and predicted positive
3. FN / False Negative: case was positive but predicted negative
4. FP / False Positive: case was negative but predicted positive

$$\text{recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

$$\text{precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

$$F_1 = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

Learnings from the assignment :

- We implemented the LSH theory which was taught in class.
- We became familiar with Pyspark programming and tried our hands on basic Pyspark functionalities and libraries available in it.

- We learned about the various preprocessing libraries and techniques that have to be applied to the data in order to get a better picture of the real content that we wish to study as there are lots of symbols, URLs and other punctuations which make it difficult to get the real content from the data.
- We initially got some errors while getting the confusion matrix
- We tried our best to improve the accuracy and get the `is_duplicate` value similar to the results as achieved by us from our own implementation of LSH which helped us explore various ways of improving the efficiency
- Initially we had difficulty in setting up the connection but we learned it about how to do it and then we uploaded our data to MongoDB and then read that data using spark into a dataframe and then we worked on this obtained dataframe.