

Locality-Sensitive Hashing

Question 1:

Here is a matrix representing the signatures of seven columns, C1 through C7.

C1	C2	C3	C4	C5	C6	C7
1	2	1	1	2	5	4
2	3	4	2	3	2	2
3	1	2	3	1	3	2
4	1	3	1	2	4	4
5	2	5	1	1	5	1
6	1	6	4	1	1	4

Suppose we use locality-sensitive hashing with three bands of two rows each. Assume there are enough buckets available that the hash function for each band can be the identity function (i.e., columns hash to the same bucket if and only if they are identical in the band). Find all the candidate pairs.

Solution:

c6 and c7

c2 and c5

c1 and c7

c3 and c7

Question 2:

Suppose we have computed signatures for a number of columns, and each signature consists of 24 integers, arranged as a column of 24 rows. There are N pairs of signatures that are 50% similar (i.e., they agree in half of the rows). There are M pairs that are 20% similar, and all other pairs (an unknown number) are 0% similar.

We can try to find 50%-similar pairs by using Locality-Sensitive Hashing (LSH), and we can do so by choosing bands of 1, 2, 3, 4, 6, 8, 12, or 24 rows. Calculate approximately, in terms of N

and M, the number of false positive and the number of false negatives, for each choice for the number of rows. Then, suppose that we assign equal cost to false positives and false negatives (an atypical assumption). Which number of rows would you choose if M:N were in each of the following ratios: 1:1, 10:1, 100:1, and 1000:1?

Question 3:

Find the set of 2-shingles for the "document":

ABRACADABRA

and also for the "document":

BRICABRAC

Answer the following questions:

1. How many 2-shingles does ABRACADABRA have?

Ans: ABRACADABRA has 5 , 2 shingles

2. How many 2-shingles does BRICABRAC have?

Ans: BRICABRAC has 3 , 2 shingles

3. How many 2-shingles do they have in common?

Ans: They hav 5 shingles in common

4. What is the Jaccard similarity between the two documents"?

Ans: Jaccard similiarity is 4/7 ands 5/7 respectively

We will use k-shingles (sometimes called n-grams) or groups of k letters

Jaccard Similarity = sets intersection / sets union Similar documents have similar k-shingles Change a word only affects k-shingles within distance k

from the word and reordering paragraphs only affects the 2k shigles that cross paragraphs boundaries

Question 4:

Consider the following matrix:

	C1	C2	C3	C4
R1	0	1	1	0
R2	1	0	1	1
R3	0	1	0	1
R4	0	0	1	0
R5	1	0	1	0
R6	0	1	0	0

Compute the Jaccard similarity between each pair of columns.

Jaccard similarity = $|a \sim b| / |a \cup b|$

Question 5: Consider the following matrix:

	C1	C2	C3	C4
R1	0	1	1	0
R2	1	0	1	1
R3	0	1	0	1
R4	0	0	1	0
R5	1	0	1	0
R6	0	1	0	0

Perform a minhashing of the data, with the order of rows: R4, R6, R1, R3, R5, R2.

Note: we give the minhash value in terms of the original name of the row, rather than the order of the row in the permutation. These two schemes are equivalent, since we only care whether hash values for two columns are equal, not what their actual values are.

Sol:

The minhash value for C1 is R4

The minhash value for C3 is R4

The minhash value for C4 is R2

The minhash value for C4 is R5