

Locality Sensitive Hashing

Rebecca C. Steorts

October 10, 2018

1 Locality Sensitive Hashing (LSH)

A *hash function* maps objects to integers such that dissimilar objects are mapped far apart. LSH uses special hash functions to ensure that similar objects are put close to each other (with high probability). By applying several locality-sensitive hashes to a record, one goes from a high-dimensional object to a low-dimensional signature, with a guarantee that similar records have nearby signatures. The low-dimensional signature vectors can then be divided into bins or blocks, with a high probability that all records mapped to the same bin are similar, and a hope that not too many similar records fall into different bins. If the number of records per bin is small, we can treat the bins as blocks for purposes of record linkage. The number of records per bin depends on the exact binning procedure and its tuning parameters (???). Unlike conventional blocking, LSH uses all the fields of a record and can be adjusted to ensure that blocks are manageably small.¹ By design, records falling within the same block are similar to each other, so only linking records within blocks can lead to dramatic speed-ups while imposing only a small cost in false negative errors. (One does fewer comparisons, and those comparisons are more likely to lead to links.)

While these are all desirable features, LSH is in several ways not yet fully developed as a blocking technique for entity resolution. One is that “similarity” must be made computationally precise, and not every similarity measure is preserved by a (known) family of hash functions. Moreover, different similarity measures may be appropriate for different kinds of data in different applications. Second, entity resolution should come not from ranking similarity scores (as in ?), but from a statistical model. We develop a

¹The last point needs some care, since simulation studies show that *sometimes* bins can contain many records, leading to minimal computational savings.

29 variety of LSH blocking algorithms which are designed to work well with real-
30 world examples of entity resolution and to meet the needs of my statistical
31 models.

32 1.1 Minwise hashing

33 Minwise hashing is a procedure for comparing the similarity of two sets
34 (records in this context can be thought of as sets). For a large dataset like
35 the Syrian one we are considering, minwise hashing is optimal because of its
36 speed in subdividing records into blocks. The implementation of minwise
37 hashing has a number of steps.

38 First

39 We shingle each record. Shingling is a way of splitting apart a record
40 into smaller subsets, at some specified splitting point. For example,
41 if we have the record [“Peter”, “Pittsburgh”], and we want to use
42 a shingling of size $k = 3$, we form the shingles [“Pet”, “erP”, “itt”,
43 “sbu”, “rgh”]. To do this, we combine the entire record into a single
44 string, and then subdivide it into parts of three.

45 Shingling is done to each record in the dataset, which forms a large set
46 of shingles. We also want to keep track of which records are associated
47 with which shingles. So if we have a second record, [“Steve”, “Pitts-
48 burgh”], then we know that the shingles “itt”, “sbu”, and “argh” are
49 in common between the two records.

50 Second

51 We form a characteristic matrix from the shingles. A characteristic
52 matrix is an indicator matrix, where the columns of the matrix corre-
53 spond to records in the dataset (so there is a column for each record)
54 and the row to the shingles formed from the records (so there is a
55 row for each shingle). The elements of the matrix are a binary (1, 0)
56 decision, where a 1 is present if the record contains the corresponding
57 shingle, and a 0 if not. The characteristic matrix is very sparse (i.e.,
58 it contains mainly 0s), as most records do not contain the majority of
59 all possible shingles.

60 Table 1 is an example of a characteristic matrix, with four records(sets)
61 and five shingles(elements).

62 Third

63 We permute the rows of the characteristic matrix to form a permuted

Element 1-5	S_1	S_2	S_3	S_4
1	0	0	1	1
2	1	0	0	1
3	1	1	0	0
4	0	0	0	1
5	1	0	0	1

Table 1: A sample characteristic matrix, with four records in columns and five possible shingles in rows. As the number of shingles grows, the matrix becomes increasingly sparse, as most records do not contain most shingles.

64 matrix. If we have five rows, 1-5, then possible permutations of the
65 rows are 12345, 15234, and 54321 (there are many others). The per-
66 muted matrix, then, is simply a reordering of the original characteristic
67 matrix, with the rows swapped in some arrangement. Figure ?? shows
68 the characteristic matrix converted to a permuted matrix by a given
69 permutation.

70 We repeat the permutation step for several iterations to obtain multi-
71 ple permuted matrices.

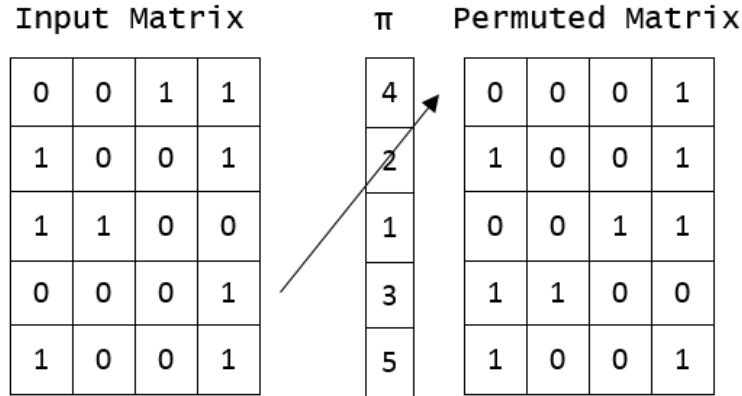


Figure 1: Permuted matrix from the characteristic one. The π vector is the specified permutation.

72 In practice, the implementation of minhashing is slightly different from
73 these permutations, as permuting the matrices a large number of times

74 is computationally-prohibitive. To make up for this, a series of hash
75 functions is generated, where the hash function assign values to each
76 row, in exactly the same way as the permutations do. The hash func-
77 tion, $(r + 1)\%5$, for example, takes in a row index (1-5), and assigns
78 the new indices as the permuted matrix. This results in the same out-
79 comes as the one above. An example of these hash functions with the
80 characteristic matrix is given in table 2.

Row	S_1	S_2	S_3	S_4	$(r+1) \% 5$	$(3r+2) \% 5$
0	0	0	1	1	1	2
1	1	0	0	1	2	0
2	1	1	0	0	3	3
3	0	0	0	1	4	1
4	1	0	0	1	0	4

Table 2: Characteristic matrix with four sets, two hash functions, and a generated binary matrix. The hash values are found by inputting the corresponding row index to the hash function.

81 **Fourth**

82 We compute the signature matrix. The signature matrix is a hashing
83 of values from the permuted one. The signature has a row for the
84 number of permutations calculated, and a column corresponding with
85 the columns of the permuted matrix. We iterate over each column of
86 the permuted matrix, and populate the signature matrix, row-wise,
87 with the row index from the first 1 value found in the column. The
88 row index inputted to the signature matrix is the new row index the
89 row was associated with, rather than the original index value. The
90 signature matrix for table 1 is in table 3.

	S_1	S_2	S_3	S_4
h_1	2	4	3	1

Table 3: Signature matrix formed from table 1. The value for each element is the row index in which the first 1 is found in the permuted matrix.

91 **Fifth**

92 Based on the results of the signature matrix, we observe the pairwise
 93 similarity of two records. For the number of hash functions, there is
 94 an interesting relationship that we use between the columns for any
 95 given set of the signature matrix and a Jaccard Similarity measure.

The Jaccard Similarity, which is a measure of similarity, is fundamental in this approach. For two sets S and T, the Jaccard Similarity is given by

$$\frac{|S \cap T|}{|S \cup T|}$$

96 which is the intersection of the two sets over their union. It provides
 97 a sense of how much two records agree in their fields, over their total
 98 number of variables.

99 The relationship between the random permutations of the character-
 100 istic matrix and the Jaccard Similarity is:

$$Pr\{min[h(A)] = min[h(B)]\} = \frac{|A \cap B|}{|A \cup B|}$$

101 The equation means that the probability that the minimum values of
 102 the given hash function, in this case h , is the same for sets A and
 103 B is equivalent to the Jaccard Similarity, especially as the number of
 104 record comparisons increases. The output of this formula, in terms of
 105 the signature matrix, can be seen in figure ??.

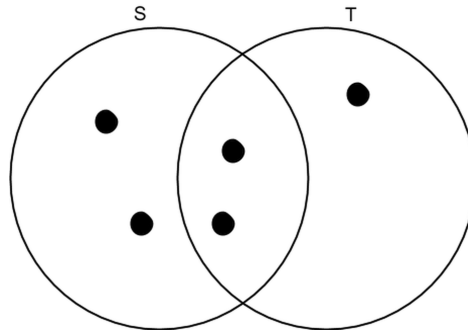


Figure 2: Two sets S and T with Jaccard similarity 2/5. The two sets share two elements in common, and there are five elements in total.

106 We use this relationship to calculate the similarity between any two
 107 records. We look down each column, and compare it to any other col-

108 umn: the number of agreements over the total number of combinations
109 is equal to Jaccard measure.

110 Of course, in practice, we would need to repeat the above for many per-
111 mutations π such that we avoid collisions, or mapping the different records
112 into the same bin.

113 Notice that we have walked through a very elementary view of locality
114 sensitive hashing, and we still have the limitation that we are having to do
115 all-to-all comparisons. How do we avoid this?

116 **2 LSH for Minhash Signatures**

117 See Section 3.4.1 in Mining for Massive Datasets and work through the
118 exercises (some of these we did in class). See Example 3.11 and Exercise
119 3.4.1 and 3.4.2.

120 **3 Further reading**

121 For further reading about hashing and it's use in practice, please refer to
122 <https://arxiv.org/abs/1710.02690>.

123 **4 Overview of Locality Sensitive Hashing (LSH)** 124 **Algorithm**

125 Below we outline the LSH algorithm (for all-to-all comparisons) and also for
126 filtering out documents that are not very similar.

- 127 1. Construct shingles of all documents in your corpus.
- 128 2. Hash all of your shingled documents.
- 129 3. Compute pairwise Jaccard similarity coefficients for all documents.
 - 130 (a) To do this in a computationally more efficient way, use the char-
131 acteristic matrix and a random permutation.
 - 132 (b) Then create the signature matrix by using the minhash. Repeat
133 this process using many random permutations in order to avoid
134 collisions. This will increase the size of your signature matrix.

135 To avoid performing all-to-all comparisons as in the above situation, we
136 compute the Jaccard similarity only for candidate pairs using b bands and
137 r rows of the signature matrix, which provide a threshold $t = (1/b)^{1/r}$ us-
138 ing the steps above but now using these extra conditions of filtering out
139 documents that are unlikely to be the same.