

similarity search reading

wbg231

January 2023

1 finding nearest neighbors

- this reading is cool, but has some unneeded section

Applications of Near-Neighbor Search

- the Jaccard similarity of sets S and T is given by

$$SIM(S, T) = \frac{|S \cap T|}{|S \cup T|}$$

that is the ratio of there intersection over there union

- this can be used to find the similarity between two documents on the characters not contextual level
- could be helpful to find plagiarism or identify bad articles from the same source
- another case where Jaccard similarity could be useful is collaborative filtering that is the task of recommending users it's based off of the items liked by other users.

shingling of documents

- k-shingles are the set of strings of length k that appear in a document
- picking the right number of shingles is critical for how useful they are
- instead of working with the shingles directly we hash them, and work at treat the key number they were hashed to as representative of that shingle
- this allows for quick and easy data compression
- could also look at words in a document as shingles

similarity preserving summary of sets

- sets of shingles are really large, even hashing them the space needed to store four byte shingles is the same as that taken to store the document

<i>Element</i>	S_1	S_2	S_3	S_4
<i>a</i>	1	0	0	1
<i>b</i>	0	0	1	0
<i>c</i>	0	1	0	1
<i>d</i>	1	0	1	1
<i>e</i>	0	0	1	0

-
- we can represent sets of characters using a characteristic or one hot encoding matrix.
- this is not actually how we store the data, but it is a good representation to keep in mind

min hashing

- to [minhash](#) a set represented as a column in a characteristic matrix, we first pick a permutation of the rows, the minhash value of any column is the number of the first row in the permuted order in which the column has 1,
- suppose we permuted the table above to be bedac the hash of $h(S_1)=a$
 $h(S_2)=c$ $h(S_3)=b$ $h(S_4)=a$
- a cool fact is that the probability that the minhash function for a random permutation of the rows produces the same value for two sets equals the Jaccard similarity of those two sets
- so in other words if we can compute the min hash of a random permutation of the characteristic matrix for the sets then we can approximate the Jaccard similarity of the sets without a lot of storage

min hash signatures

- so we can represent the characteristic matrix of a set of vectors by picking n random permutations of the M rows

- then we minhash these permutations, then for the col represent S we construct the min hash signature for S as the vector of all the minhashes of s
- thus we can think of the matrix m as the matrix with all signature vectors in its columns
- notice this compresses our representation matrix
- we can not do the minhash in practice
- but we can simulate the effect of a random permutation by a random hash function that maps row number to buckets
- so this hash function approximates a permutation
- so instead of picking n random permutations we pick n random hash functions on the rows.
- we construct the signature matrix by considering each row in their given order.
- we handle row r by computing all hashes for row r then for each col if c has a 0 in row r do nothing, if c has a 1 in row r set $\text{sig}(i,c)$ to be the min of the current value and the hash of every row at h
- i am a bit confused about the computation of the min hashing

local sensitive hashing for documents

- in cases with really large number of documents or items it may be infeasible to even compare the signature matrix
- this is where Local sensitive hashing comes in
- we can just get the most similar items using local sensitive hashing

LSH for minihash signatures

- a general approach to lsh is to hash items several times in such a way that similar items are more likely to be hashed to the same bucket than dissimilar items
- then we just check the similarity of pairs that hash to the same bucket

distance metrics

- there are a number of other worth while distance metric
- l2 distance
- jacard distance (1-jackard similarity)
- cosine distance which i proprtional to teh angle between two vectors
- edit distance is how many insertions or deletions must be made to get from one string ot the other
- hamming distance number of components in which two vectors differ
-