

Click here [Slide\(pdf\)](#))

LSH: Locality-Sensitive Hashing

a way to find "similar" sets.

upside: only a small fraction of points are ever examined

downside: exists false negatives

Steps for similar docs

first. Shingling

Convert a document into a set

1. hash shingles to a few bytes

2. Compute $\text{Sim}(C1, C2)$

How to compute?

First define the Jaccard Similarity.

$\text{JS}(\text{not dis}) = \text{num of (equal items which not 0)} / \text{num of (all items)}$

$\text{dis} = 1 - \text{JS}.$

The smaller the dis, the more similar the two vectors are.

second. Min-hashing

Convert large sets to short signatures, while preserving similarity.

Find a hash function $h(x)$ make that:

if $\text{sim}(C1, C2)$ is high, with high prob. $h(c1)$ equal $h(C2)$;

if low, with high prob. $h(c1) \neq h(c2)$

The function to our taste is Min-hashing

Then we get signatures.

third. LSH :

Focus on pairs of signatures likely to be from similar documents

There uses bands and rows to reduce the error.

See the slides page 40.

Use hashing to find candidate pairs of similarity $\geq s$