

# INFORMATION RETRIEVAL

## Near duplicate detection

---

Corso di Laurea Magistrale in Informatica

Università di Roma Tor Vergata

Prof. Giorgio Gambosi

a.a. 2021-2022

Derived from slides originally produced by C. Manning and by H. Schütze



Many problems in **data mining** can be seen as searching in sets of **similar** items:

- ⊙ Pages with similar words, for **classification** on topics.
- ⊙ Topic suggestion to Twitter users with similar profiles (**recommendation systems**).
- ⊙ Dual problem: identifying **communities** of users with similar interests
- ⊙ Identifying same user in different contexts (e.g. social media platforms)

- ⊙ The web is full of duplicated content.
- ⊙ More so than many other collections
- ⊙ Exact duplicates
  - Easy to eliminate
  - E.g., use hash/fingerprint
- ⊙ Near-duplicates
  - Abundant on the web
  - Difficult to eliminate
- ⊙ For the user, it's annoying to get a search result with near-identical documents.
- ⊙ **Marginal relevance is zero**: even a highly relevant document becomes nonrelevant if it appears below a (near-)duplicate.
- ⊙ We need to eliminate near-duplicates.

Finding sets of documents (web pages) with much text in common:

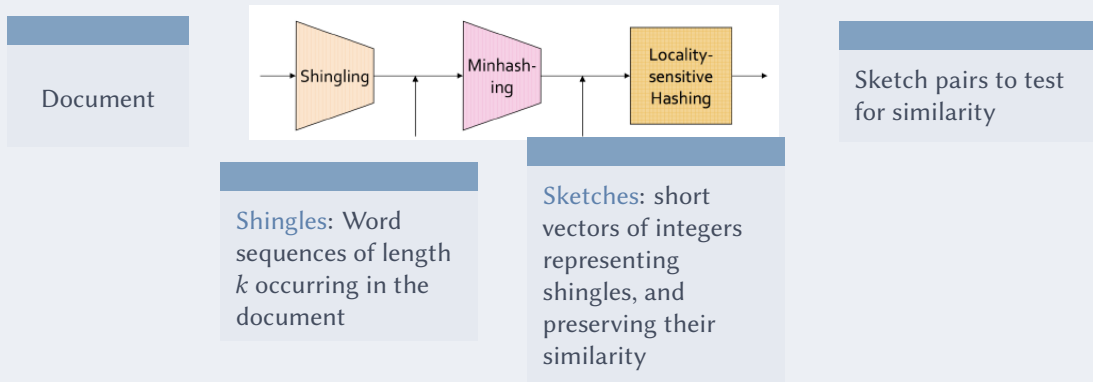
- ⊙ Mirror or quasi-mirror sites
  - **Application:** elimination of duplicates.
- ⊙ **Plagiarism**, inclusion of extensive citations .
- ⊙ Articles with similar content in different news sites .
  - **Application:** grouping articles as a “common history”.

# Detecting near-duplicates

- ⊙ Compute similarity with an edit-distance measure
- ⊙ We want “syntactic” (as opposed to semantic) similarity.
  - True semantic similarity (similarity in content) is too difficult to compute.
- ⊙ We do not consider documents near-duplicates if they have the same content, but express it with different words.
- ⊙ Use similarity threshold  $\theta$  to make the call “is/isn’t a near-duplicate”.
- ⊙ E.g., two documents are near-duplicates if similarity  $> \theta = 80\%$ .

## Three techniques useful for NDD

- ⊙ **Shingling**: convert documents, e-mail, ecc, in sets of items.
- ⊙ **Minhashing**: convert large sets in short **sketches** (or signatures), preserving similarity.
- ⊙ **Locality Sensitive Hashing (LSH)**: consider pairs of signature that could be similar with at least a given probability.



# Represent each document as set of shingles

Shingles are used as features to measure **syntactic similarity** of documents.

- ⊙ A shingle is just a **word  $k$ -gram**.
- ⊙ A document is represented as a set of shingles
- ⊙ For  $n = 5$ , “*In a hole in the ground there lived a hobbit*” would be represented as this set of shingles:
  - {In a hole in the, a hole in the ground, hole in the ground there, in the ground there lived, the ground there lived a, ground there lived a hobbit }
- ⊙ Similar documents will have many shingles in common



## Represent each document as set of shingles

- ⊙ Modifying a word affects only  $k$  shingles (the ones at distance at most  $k$  from the word)
- ⊙ Moving a paragraph affects  $2k$  shingles (the ones at distance at most  $k$  from the paragraph borders)
- ⊙ For  $n = 3$ , changing “*In a hole in the ground there lived a hobbit*” to “*In a hole in the ground there was a hobbit*” only changes shingles { ground there lived, there lived a, lived a hobbit }

- ⊙ In general, different documents should have few shingles in common, especially for higher  $k$
- ⊙ We define the similarity of two documents as the **Jaccard** coefficient of their shingle sets.

## Recall: Jaccard coefficient

- ⊙ A commonly used measure of overlap of two sets
- ⊙ Let  $A$  and  $B$  be two sets: their Jaccard coefficient is defined as:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$(A \neq \emptyset \text{ or } B \neq \emptyset)$

- ⊙  $J(A, A) = 1$
- ⊙  $J(A, B) = 0$  if  $A \cap B = \emptyset$
- ⊙  $A$  and  $B$  don't have to be the same size.
- ⊙ Always assigns a number between 0 and 1.

# Jaccard coefficient: Example

- ⊙ Three documents:
  - $d_1$ : “Jack London traveled to Oakland”
  - $d_2$ : “Jack London traveled to the city of Oakland”
  - $d_3$ : “Jack traveled from Oakland to London”
- ⊙ Based on shingles of size 2 (2-grams or bigrams), what are the Jaccard coefficients  $J(d_1, d_2)$  and  $J(d_1, d_3)$ ?
- ⊙
  1.  $s(d_1) = \{\text{“Jack London”, “London traveled”, “traveled to”, “to Oakland”}\}$
  2.  $s(d_2) = \{\text{“Jack London”, “London traveled”, “traveled to”, “to the”, “the city”, “city of”, “of Oakland”}\}$
  3.  $s(d_3) = \{\text{“Jack traveled”, “traveled from”, “from Oakland”, “Oakland to”, “to London”}\}$
- ⊙ there are
- ⊙  $J(d_1, d_2) = 3/8 = 0.375$
- ⊙  $J(d_1, d_3) = J(d_2, d_3) = 0$

## Represent each document as a sketch

- ⊙ The number of shingles per document is large: computing Jaccard directly from  $M$  is expensive
- ⊙ To increase efficiency, we will represent documents by means of **sketches**, cleverly chosen **subsets** of their shingles.
- ⊙ Let  $h$  be a predefined sketch size and let  $S$  be the overall set of shingles: document sketches are derived by means of a set of  $h$  different **random permutations**  $\pi_1 \dots \pi_h$  of  $S$
- ⊙ Each  $\pi_i$  maps a shingle to a different integer in  $\{1, \dots, |S|\}$
- ⊙ The **sketch** of a document  $d$  is defined as:

$$\left( \min_{s \in d} \pi_1(s), \min_{s \in d} \pi_2(s), \dots, \min_{s \in d} \pi_h(s) \right)$$

(a vector of  $h$  integers).

## From sets of documents+shingles to boolean matrices

A set of documents can be represented as a boolean matrix  $M$ , where

- ⊙ columns are associated to documents
- ⊙ rows correspond to all shingles appearing in any document
- ⊙  $M(i, j) = 1$  iff the  $i$ -th shingle appear in the  $j$ -th document
- ⊙ The matrix is usually sparse

The Jaccard **similarity** of two documents can be derived from the corresponding columns

## Four types of rows

- For any pair of columns  $S_1, S_2$ , rows can be classified in four types according to the values of the corresponding values in the matrix: each type has a different effect on numerator  $N$  and denominator  $D$  of  $J(S_1, S_2)$

	$S_1$	$S_2$	effect on $N$	effect on $D$
a	1	1	increase	increase
b	1	0	same	increase
c	0	1	same	increase
d	0	0	same	same

- In fact,  $J(S_1, S_2) = \frac{\#a}{\#a + \#b + \#c}$
- Many rows are of type  $d$

Permutations of shingles correspond here to permutations of rows of  $M$ . The above considerations can be accordingly translated as follows.

- ⊙ Given a row permutation  $\pi$ , for any document  $d$  corresponding to a column  $c_i$  in  $M$ , let us define as the **Minhash** of  $d$  under permutation  $\pi$ , denoted as  $MH_\pi(d)$  the index  $j$  of the **first** row (according to  $\pi$ ) such that  $M(j, i) = 1$ .
- ⊙ As an extension, given a set  $\Pi_r$  of  $r$  permutations, for any document  $d$  corresponding to a column  $c_i$  in  $M$ ,  $MH_{\Pi_r}(d)$  is defined as the vector of integers  $(j_1, \dots, j_r)$  such that  $j_t$  is the index of the **first** row (according to permutation  $\pi_t$ ) such that  $M(j_t, i) = 1$ .



- ⊙ The sketch vector  $MH_{\Pi_r}(d)$  can be interpreted as a signature of  $d$
- ⊙ Signatures can be visualized as columns in a new matrix  $M'$ , where columns correspond to documents while rows correspond to permutations. The values in column  $c_i$  are then defined as  $MH_{\Pi_r}(d_i)$ , where  $d_i$  is the document corresponding to  $c_i$

Shingle/document  
matrix  $M$

$d_1$	$d_2$	$d_3$	$d_4$
1	0	1	0
1	0	0	1
0	1	0	1
0	1	0	1
0	1	0	1
1	0	1	0
1	0	1	0

# Minhashing example

Permutations

1
3
7
6
2
5
4

$M$

$d_1$	$d_2$	$d_3$	$d_4$
1	0	1	0
1	0	0	1
0	1	0	1
0	1	0	1
0	1	0	1
1	0	1	0
1	0	1	0

Signature matrix  $M'$

$s_1$     $s_2$     $s_3$     $s_4$

1	2	1	2
---	---	---	---

# Minhashing example

Permutations

1	4
3	2
7	1
6	3
2	6
5	7
4	5

$M$

$d_1$	$d_2$	$d_3$	$d_4$
1	0	1	0
1	0	0	1
0	1	0	1
0	1	0	1
0	1	0	1
1	0	1	0
1	0	1	0

Signature matrix  $M'$

$S_1$     $S_2$     $S_3$     $S_4$

2	1	4	1
1	2	1	2

# Minhashing example

Permutations

1	4	3
3	2	4
7	1	7
6	3	6
2	6	1
5	7	2
4	5	5

$M$

$d_1$	$d_2$	$d_3$	$d_4$
1	0	1	0
1	0	0	1
0	1	0	1
0	1	0	1
0	1	0	1
1	0	1	0
1	0	1	0

Signature matrix  $M'$

$S_1$	$S_2$	$S_3$	$S_4$
2	1	2	1
2	1	4	1
1	2	1	2

# Detecting near-duplicates from sketches

Assume a single permutation  $\pi$ . Check is performed as follows:

- ⊙ If  $\text{MH}_{\pi(d_1)} = \text{MH}_{\pi(d_2)}$  then  $d_1$  and  $d_2$  probably are near-duplicates.
- ⊙ If  $\text{MH}_{\pi(d_1)} \neq \text{MH}_{\pi(d_2)}$  then  $d_1$  and  $d_2$  are probably not near-duplicates.

# Detecting near-duplicates from sketches

Why does it work? Let us first recall that by  $b, c, a$  we denote the set of shingles in  $d_1$  and not in  $d_2$ , in  $d_2$  and not in  $d_1$ , in both  $d_1$  and  $d_2$ , respectively. Then,

- ⊙ the number of shingles occurring in  $d_1$ , that is of rows  $i$  such that  $M'(i, 1) = 1$ , is  $\#a + \#b$
- ⊙ similarly, the number of shingles occurring in  $d_2$ , that is of rows  $i$  such that  $M'(i, 2) = 1$ , is  $\#a + \#c$
- ⊙ the number of possible (not distinct) pairs of shingles, the first one occurring in  $d_1$  and the second one in  $d_2$ , that is of (not distinct) pairs of rows  $i, j$  in  $M'$  such that  $M'(i, 1) = M'(j, 2) = 1$  is  $(\#a + \#b)(\#a + \#c) - \#b\#c$
- ⊙ the number of possible (not distinct) pairs of shingles both occurring in both  $d_1$  and in  $d_2$ , that is of (not distinct) pairs of rows  $i, j$  in  $M'$  such that  $M'(i, 1) = M'(i, 2) = M'(j, 1) = M'(j, 2) = 1$  is  $\#a^2$

# Detecting near-duplicates from sketches

Let us now estimate the probability that, by randomly choosing  $\pi$ , we get  $MH_{\pi}(d_1) = MH_{\pi}(d_2)$ .

- ⊙ the number of possible pairs  $(MH_{\pi}(d_1), MH_{\pi}(d_2))$  is equal to the number of pairs of shingles, the first one occurring in  $d_1$  and the second one in  $d_2$ , that is  $(\#a + \#b)(\#a + \#c) - \#b\#c$
- ⊙ the number of possible pairs  $(MH_{\pi}(d_1), MH_{\pi}(d_2))$  with  $MH_{\pi}(d_1) = MH_{\pi}(d_2)$  is equal to the number of pairs of shingles both occurring in both  $d_1$  and in  $d_2$ , that is  $\#a^2$
- ⊙ assuming a uniform probability of selection of permutations, the probability that  $MH_{\pi}(d_1) = MH_{\pi}(d_2)$  is then given by

$$p(d_1, d_2) = \frac{\#a^2}{(\#a + \#b)(\#a + \#c) - \#b\#c}$$



- ⊙ But

$$\frac{\#a^2}{(\#a + \#b)(\#a + \#c) - \#b\#c} = \frac{\#a}{\#a + \#b + \#c}$$

is the Jaccard coefficient  $J(d_1, d_2)$ , that is our similarity measure between  $d_1$  and  $d_2$ .  
So, estimating  $p(d_1, d_2)$  corresponds to estimating the similarity between  $d_1$  and  $d_2$

- ⊙ How can we get a good estimate of  $p(d_1, d_2)$  more efficiently than computing  $J(d_1, d_2)$  (which implies taking into account all their shingles?)

## Detecting near-duplicates from sketches

- ⊙  $p(d_1, d_2)$  can be seen as the probability that, given  $d_1$  and  $d_2$ , a uniformly sampled permutation of the set of shingles assigns the same index to the first shingle in both documents
- ⊙ Selecting  $\pi$  and observing whether  $MH_\pi(d_1) = MH_\pi(d_2)$  can be seen as sampling a stone from an urn containing  $\#a$  red stones and  $\#b + \#c$  black stones and checking whether the sampled stone is red

## Detecting near-duplicates from sketches

- ⊙ Performing a random sample of  $r$  independent permutations  $\pi_1, \dots, \pi_r$  and observing whether  $MH_{\pi_i}(d_1) = MH_{\pi_i}(d_2)$  for each  $\pi_i$  corresponds to sampling  $r$  stones from the urn (with replacement) and checking how many sampled stones are red
- ⊙ This is a sequence of **Bernoulli trials** with probability  $p(d_1, d_2)$ . In this case, the number of red stones is distributed according to a binomial distribution

$$p(MH_{\pi_i}(d_1) = MH_{\pi_i}(d_2) \text{ for } t \text{ permutations}) = \binom{t}{r} p(d_1, d_2)^t (1 - p(d_1, d_2))^{r-t}$$

which has mean  $r \cdot p(d_1, d_2)$

## Detecting near-duplicates from sketches

- ⊙  $J(d_1, d_2)$  can be estimated by estimating  $p(d_1, d_2)$  from the sample of size  $r$  provided by the set functions  $\Pi_r$ .
- ⊙ by standard statistics, an unbiased estimator of  $p$  is  $\hat{p} = \frac{t}{r}$ , where  $t$  is the number of functions  $h \in \Pi_k$  such that  $\text{MH}_\pi(d_1) = \text{MH}_\pi(d_2)$
- ⊙ the corresponding standard error is given by the sample standard deviation  $\hat{s} = \sqrt{\frac{\hat{p}(1-\hat{p})}{r}}$ : this makes it possible to define a confidence interval on  $J(d_1, d_2)$  at any given confidence level  $\theta$  as  $[\hat{p} - Z_\theta \hat{s}, \hat{p} + Z_\theta \hat{s}]$ , where  $Z_\theta$  is the  $Z$ -score at probability  $\theta$  (number of standard deviations from the mean of a gaussian such that the tail probability is  $1 - \theta$ )
- ⊙ the precision of the estimation improves as  $r$  increases

# Random hash functions as permutations

Sketches can be efficiently computed by means of **random hash functions**.

- ⊙ We can map shingles to integers by fingerprinting, that is by applying a given hash function  $h$  which maps any sequence of unigrams to a sequence of (say)  $m$  bytes, that is to an integer interval  $0..2^m - 1$
- ⊙ For suitably large  $m$ , with high probability there is no collision between pairs of shingles, that is  $h(s_1) \neq h(s_2)$  for all  $s_1, s_2$
- ⊙ Then, for suitably large  $m$ ,  $h$  defines a permutation of shingles with high probability

# Implementing Minhashing

- ⊙ Let  $h_1, \dots, h_k$  be a set of **hash** functions and  $d_1, \dots, d_m$  a set of documents
- ⊙ Let  $S$  be a  $k \times m$  matrix: at the end of the algorithm: let  $S(i, j) = \infty$  for all  $i, j$
- ⊙ For each document  $d_j$ 
  - For each shingle  $s$  in  $d_j$ 
    - For each hash function  $h_i$ , set  $S(i, j) = \min(h_i(s), S(i, j))$
- ⊙ At the end,  $S(i, j)$  will store the minimum index, in the permutation of shingles induced by  $h_i$ , of a shingle in document  $d_j$ . This is the MinHash of document  $d_j$  when function  $h_i$  is applied.

## Example

i	$d_1$	$d_2$
0	1	1
1	0	0
2	1	1
3	1	0
4	0	1

$$h_1(x) = x \bmod 5$$

$$h_2(x) = (2x + 1) \bmod 5$$

$h_1$	$d_1$	$d_2$
0	1	1
1	0	0
2	1	1
3	1	0
4	0	1

$h_2$	$d_1$	$d_2$
0	0	0
1	1	0
2	1	1
3	1	1
4	0	1

$$\min(h_1(d_1))=0=0=\min(h_1(d_2))$$

$$\min(h_2(d_1))=1 \neq 2 = \min(h_2(d_2))$$

$$\hat{J}(d_1, d_2) = \frac{1}{2} = .5$$

$$J(d_1, d_2) = \frac{2}{5} = .4$$

# Example

i	$d_1$	$d_2$
0	1	1
1	0	0
2	1	1
3	1	0
4	0	1

	$M(h_i(s),1)$	$M(h_i(s),2)$	$S$	
$h_1$			$\infty$	$\infty$
$h_2$			$\infty$	$\infty$
$h_1(0) = 0$	1	1	0	0
$h_2(0) = 1$	0	0	$\infty$	$\infty$
$h_1(1) = 1$	0	0	0	0
$h_2(1) = 3$	1	0	1	$\infty$
$h_1(2) = 2$	1	1	0	0
$h_2(2) = 0$	1	1	1	2
$h_1(3) = 3$	1	0	0	0
$h_2(3) = 2$	1	1	1	2
$h_1(4) = 4$	0	1	0	0
$h_2(4) = 4$	0	1	1	2

$$h_1(x) = x \bmod 5$$

$$h_2(x) = (2x + 1) \bmod 5$$



# Efficient near-duplicate detection

- ⊙ We have an extremely efficient method for estimating similarity for a **single** pair of documents
- ⊙ But we still have to estimate  $O(n^2)$  values where  $n$  is the number of documents: still intractable
- ⊙ However, often we need to derive all pairs whose similarity is above a given threshold
- ⊙ One solution: **locality sensitive hashing** (LSH)

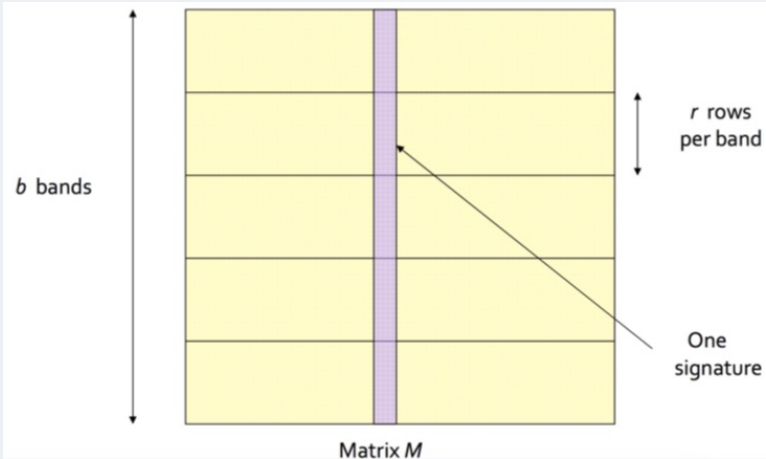
# Candidate pairs

- ⊙ pick a similarity threshold  $\theta$ ,  $0 \leq \theta \leq 1$
- ⊙ goal: find pairs of documents with Jaccard similarity at least  $\theta$
- ⊙ columns  $i$  and  $j$  are a **candidate pair** if their signatures agree in at least a fraction  $\theta$  of their rows
- ⊙ we expect pairs of documents to have the same similarity as their signatures

# Locality-Sensitive Hashing (LSH) for signatures

- ⊙ **Idea:** Hash columns of signatures matrix  $M$  to a predefined set of **buckets** in such a way that similar columns are likely to be hashed to the same bucket, with high probability
- ⊙ A pair of columns hashed to the same bucket is a **candidate pair** for similarity, to be verified more accurately
- ⊙ False positives (dissimilar pairs hashed to same bucket); false negatives (similar pairs hashed to different buckets)

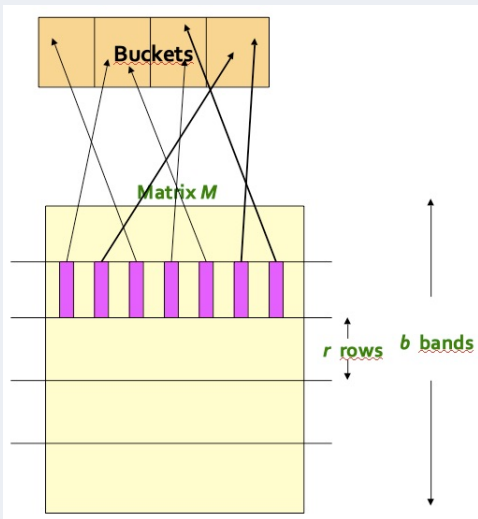
# Partition in bands



## Partition in bands

- ⊙ Divide the signature matrix  $S$  into  $b$  bands, each of  $r$  rows.
- ⊙ For each band  $B_i$ , a hash function  $h_i$  is defined which maps vectors of  $r$  integers to  $k$  buckets, with  $k$  large enough
- ⊙ We could use the same hash functions for all bands, but different bucket arrays
- ⊙ A pair of columns is a candidate pair if they are hashed to the same bucket for at least 1 band
- ⊙ Tune  $b$  (and correspondingly  $r$ ) to catch most similar pairs, but few not similar ones.

# Band hashing



- ⊙ Columns 2 and 6 are probably identical (candidate pair)
- ⊙ Columns 6 and 7 are different (wrt to this band, they could be declared candidate pairs by hashing the other bands)

## Example

- ⊙ Assume we have  $10^5$  columns (documents).
- ⊙ Each signature is a vector of length 100 (100 hash functions applied).
- ⊙ Each signature element is an integer 4 bytes long.
- ⊙ Then all signatures are 40MB long.
- ⊙ The naive approach requires  $10^5 \times (10^5 - 1) \times .5 \simeq 5 \times 10^9$  pairs of signatures to be compared: could take months
- ⊙ Let us apply LSH: choose, for example,  $b = 20, r = 5$

Assume we wish all document pairs with similarity at least .8

- ⊙ Let columns  $C_1, C_2$  be signatures of similar documents: that is, they have equal values in at least a .8 fraction of their rows
- ⊙ The probability that columns  $C_1, C_2$  collide in a given band is then  $(0.8)^5 = 0.328$ .
- ⊙ The probability that  $C_1, C_2$  do not collide in any of the 20 bands is then  $(1 - 0.328)^{20} \simeq 0.00035$ .
  - that is, there is a chance of 1 over about 3000 that two 0.8 similar columns do not collide anywhere, and are declared not similar (false negative)
  - we would find 99.965% pairs of truly similar documents: very few false negatives



# False positives

- ⊙ Assume columns  $C_1, C_2$  are signatures of not similar documents: they have equal values in a .3 fraction of their rows
- ⊙ The probability that columns  $C_1, C_2$  collide in a given band is then  $(0.3)^5 = 0.00243$ .
- ⊙ The probability that  $C_1, C_2$  collide in at least one of the 20 bands is then  $1 - (1 - 0.00243)^{20} \simeq 0.0474$ .
  - that is, approximately 4.74% pairs of docs with similarity 0.3% end up becoming candidate pairs (false positive)
  - they will be checked more precisely and it will turn out they are not similar (at .8 threshold)

## Collision probability in a band

- ⊙ The probability that two given columns  $C_1, C_2$  have equal rows in a certain band is  $\theta^r$
- ⊙ The probability that two given columns  $C_1, C_2$  differ in at least one row in a certain band is  $1 - \theta^r$
- ⊙ The probability that two given columns  $C_1, C_2$  differ in at least one row in all bands is  $(1 - \theta^r)^b$
- ⊙ The probability that two given columns  $C_1, C_2$  have equal rows in at least one band (they are a candidate pair) is  $1 - (1 - \theta^r)^b$

# LSH Involves a Tradeoff

- ⊙ Pick
  - The number of MinHashes (rows of  $S$ )
  - The number of bands  $b$
  - The number of rows  $r$  per band
- ⊙ to balance false positives/negatives
- ⊙ Example: If we had only 15 bands of 5 rows, the number of false positives would go down, but the number of false negatives would go up

**Example:**  $b = 20, r = 5$

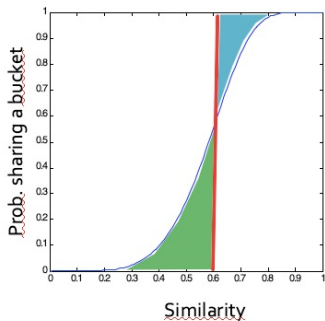
⊙ Similarity threshold  $\theta$

$\theta$	$1 - (1 - \theta^r)^b$
.2	.006
.3	.047
.4	.186
.5	.47
.6	.802
.7	.975
.8	.9996

⊙ Probability that at least 1 band is identical (collision)

# Picking the S-curve

- ⊙ Picking  $r$  and  $b$  to get the best curve
- ⊙ 50 hash-functions ( $r = 5, b = 10$ )



- Blue area: False Negative rate
- Green area: False Positive rate

- ⊙ Tune  $S, b, r$  to get almost all pairs with similar signatures, but eliminate most pairs that do not have similar signatures
- ⊙ Check in main memory that candidate pairs really do have similar signatures