

# Finding Similar Sets

Applications

Shingling

Minhashing

Locality-Sensitive Hashing

Mining of Massive Datasets

Leskovec, Rajaraman, and Ullman

Stanford University



# Applications of Set-Similarity

Many data-mining problems can be expressed as finding “similar” sets:

1. Pages with similar words, e.g., for classification by topic.
2. Netflix users with similar tastes in movies, for recommendation systems.
3. **Dual**: movies with similar sets of fans.
4. Entity resolution.

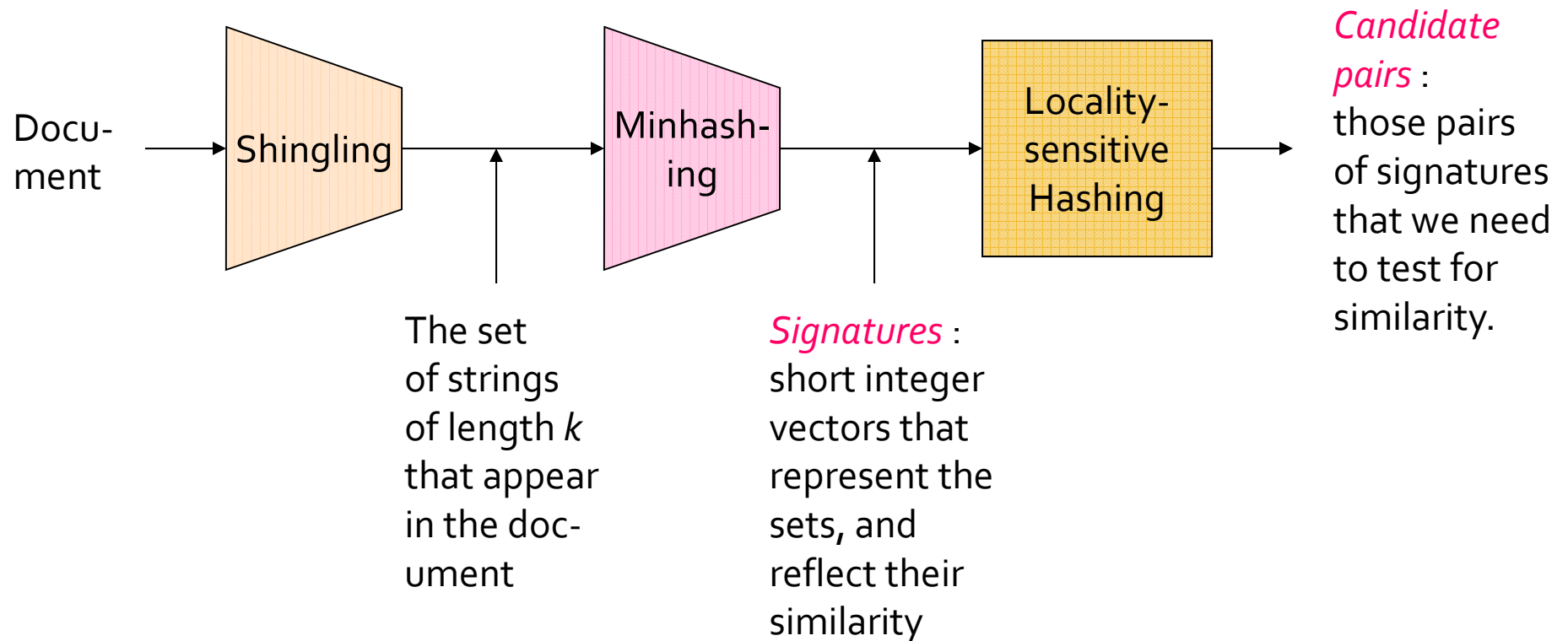
# Similar Documents

- Given a body of documents, e.g., the Web, find pairs of documents with a lot of text in common, such as:
  - Mirror sites, or approximate mirrors.
    - **Application**: Don't want to show both in a search.
  - Plagiarism, including large quotations.
  - Similar news articles at many news sites.
    - **Application**: Cluster articles by “same story.”

# Three Essential Techniques for Similar Documents

1. *Shingling* : convert documents, emails, etc., to sets.
2. *Minhashing* : convert large sets to short signatures, while preserving similarity.
3. *Locality-sensitive hashing* : focus on pairs of signatures likely to be similar.

# The Big Picture



# Shingles

- A  $k$ -shingle (or  $k$ -gram) for a document is a sequence of  $k$  characters that appears in the document.
- **Example:**  $k=2$ ; doc = abcab. Set of 2-shingles = {ab, bc, ca}.
- Represent a doc by its set of  $k$ -shingles.

# Shingles and Similarity

- Documents that are intuitively similar will have many shingles in common.
- Changing a word only affects  $k$ -shingles within distance  $k$  from the word.
- Reordering paragraphs only affects the  $2k$  shingles that cross paragraph boundaries.
- **Example:**  $k=3$ , “The dog which chased the cat” versus “The dog that chased the cat”.
  - Only 3-shingles replaced are  $g\_w$ ,  $\_wh$ ,  $whi$ ,  $hic$ ,  $ich$ ,  $ch\_$ , and  $h\_c$ .

# Shingles: Compression Option

- To compress long shingles, we can hash them to (say) 4 bytes.
  - Called *tokens*.
- Represent a doc by its tokens, that is, the set of hash values of its  $k$ -shingles.
- Two documents could (rarely) appear to have shingles in common, when in fact only the hash-values were shared.



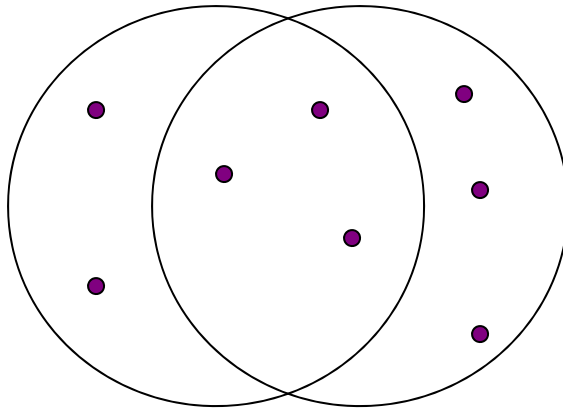
# Minhashing

Jaccard Similarity Measure  
Constructing Signatures

# Jaccard Similarity

- The *Jaccard similarity* of two sets is the size of their intersection divided by the size of their union.
- $Sim(C_1, C_2) = |C_1 \cap C_2| / |C_1 \cup C_2|.$

# Example: Jaccard Similarity



3 in intersection.  
8 in union.  
Jaccard similarity  
=  $3/8$

# From Sets to Boolean Matrices

- **Rows** = elements of the universal set.
  - **Example**: the set of all  $k$ -shingles.
- **Columns** = sets.
- 1 in row  $e$  and column  $S$  if and only if  $e$  is a member of  $S$ .
- Column similarity is the Jaccard similarity of the sets of their rows with 1.
- Typical matrix is sparse.

# Example: Column Similarity

C<sub>1</sub> C<sub>2</sub>

0 1 \*

1 0 \*

1 1 \* \*

0 0

1 1 \* \*

0 1 \*

$$\text{Sim}(C_1, C_2) = \frac{2}{5} = 0.4$$

# Four Types of Rows

- Given columns  $C_1$  and  $C_2$ , rows may be classified as:

	<u><math>C_1</math></u>	<u><math>C_2</math></u>
$a$	1	1
$b$	1	0
$c$	0	1
$d$	0	0

- Also,  $a$  = # rows of type  $a$  , etc.
- Note  $Sim(C_1, C_2) = a/(a + b + c)$  .

# Minhashing

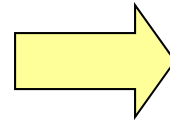
- Imagine the rows permuted randomly.
- Define *minhash function*  $h(C)$  = the number of the first (in the permuted order) row in which column  $C$  has 1.
- Use several (e.g., 100) independent hash functions to create a signature for each column.
- The signatures can be displayed in another matrix – the *signature matrix* – whose columns represent the sets and the rows represent the minhash values, in order for that column.

# Minhashing Example

Input matrix

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

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  $M$

2	1	2	1
2	1	4	1
1	2	1	2



# Surprising Property

- The probability (over all permutations of the rows) that  $h(C_1) = h(C_2)$  is the same as  $\text{Sim}(C_1, C_2)$ .
- Both are  $a / (a + b + c)!$
- Why?
  - Look down the permuted columns  $C_1$  and  $C_2$  until we see a 1.
  - If it's a type- $a$  row, then  $h(C_1) = h(C_2)$ . If a type- $b$  or type- $c$  row, then not.

# Similarity for Signatures

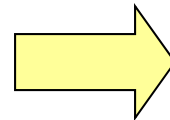
- The *similarity of signatures* is the fraction of the minhash functions in which they agree.
  - Thinking of signatures as columns of integers, the similarity of signatures is the fraction of rows in which they agree.
- Thus, the expected similarity of two signatures equals the Jaccard similarity of the columns or sets that the signatures represent.
  - And the longer the signatures, the smaller will be the expected error.

# Min Hashing – Example

Input matrix

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

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  $M$

2	1	2	1
2	1	4	1
1	2	1	2

	1-3	2-4	1-2
Col/Col	0.75	0.75	0
Sig/Sig	0.67	1.00	0

# Locality-Sensitive Hashing

Focusing on Similar Minhash Signatures  
Other Applications Will Follow

# Locality-Sensitive Hashing

- **General idea:** Generate from the collection of all elements (signatures in our example) a small list of *candidate pairs*: pairs of elements whose similarity must be evaluated.
- **For signature matrices:** Hash columns to many buckets, and make elements of the same bucket candidate pairs.

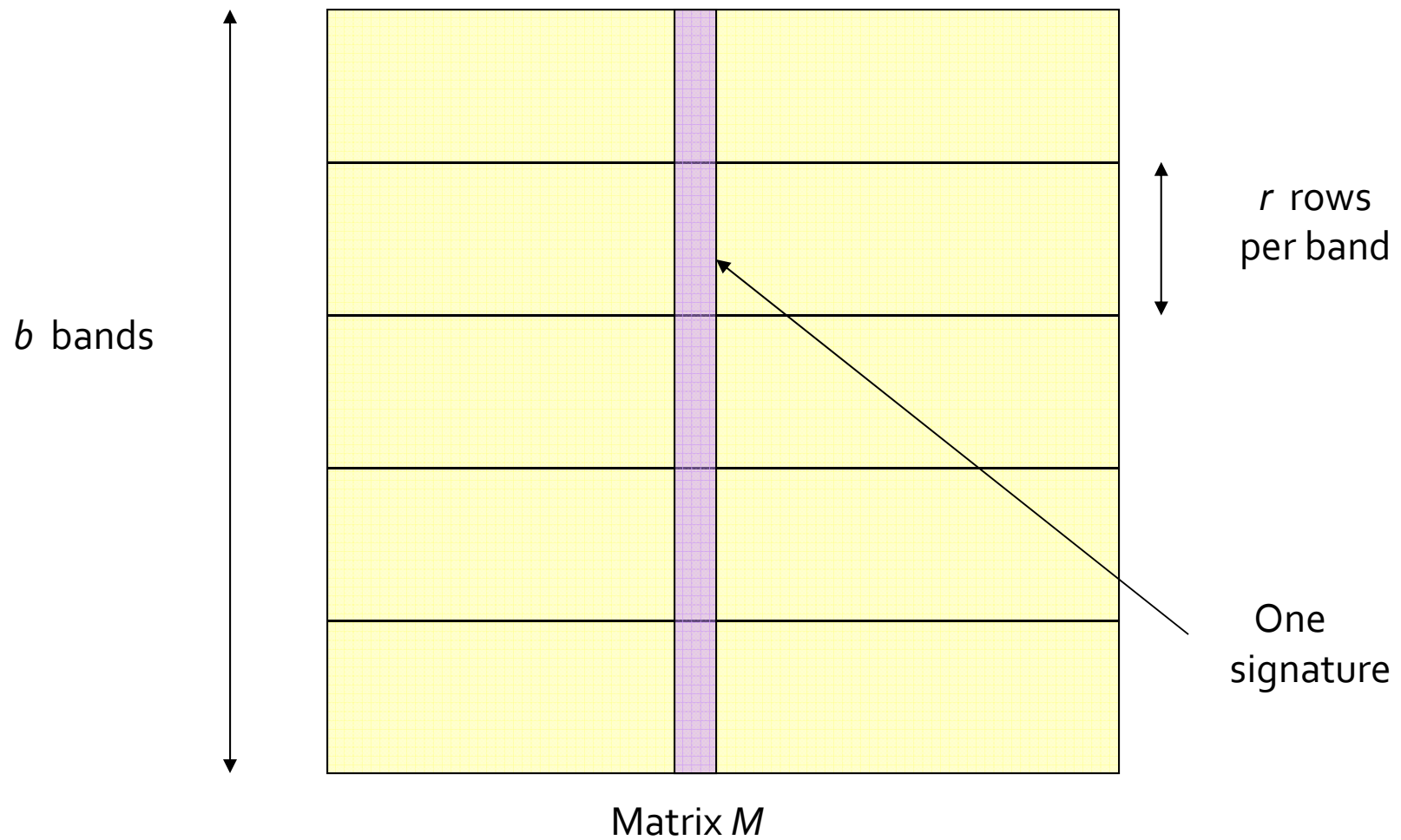
# Candidate Generation From Minhash Signatures

- Pick a similarity threshold  $t$ , a fraction  $< 1$ .
- We want a pair of columns  $c$  and  $d$  of the signature matrix  $M$  to be a *candidate pair* if and only if their signatures agree in at least fraction  $t$  of the rows.
  - I.e.,  $M(i, c) = M(i, d)$  for at least fraction  $t$  values of  $i$ .

# LSH for Minhash Signatures

- **Big idea**: hash columns of signature matrix  $M$  several times.
- Arrange that (only) similar columns are likely to hash to the same bucket.
- Candidate pairs are those that hash **at least once** to the same bucket.

# Partition Into Bands

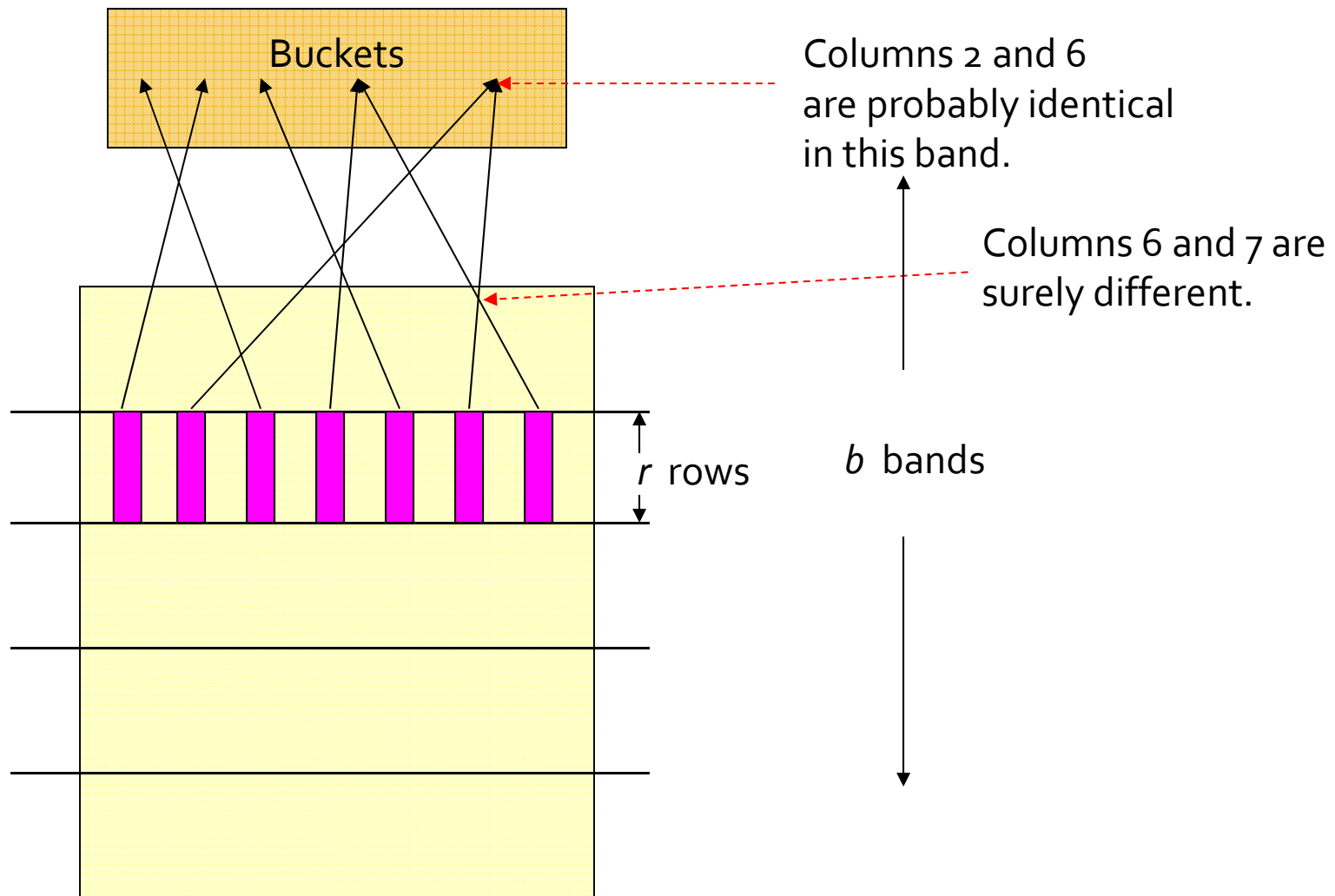




# Partition into Bands – (2)

- Divide matrix  $M$  into  $b$  bands of  $r$  rows.
- For each band, hash its portion of each column to a hash table with  $k$  buckets.
  - Make  $k$  as large as possible.
- *Candidate* column pairs are those that hash to the same bucket for  $\geq 1$  band.
- Tune  $b$  and  $r$  to catch most similar pairs, but few nonsimilar pairs.

# Hash Function for One Bucket



Matrix M

# Example – Bands

- Suppose 100,000 columns.
- Signatures of 100 integers.
- Therefore, signatures take 40Mb.
- Want all 80%-similar pairs of documents.
- 5,000,000,000 pairs of signatures can take a while to compare.
- Choose 20 bands of 5 integers/band.

# Suppose $C_1, C_2$ are 80% Similar

- Probability  $C_1, C_2$  identical in one particular band:  $(0.8)^5 = 0.328$ .
- Probability  $C_1, C_2$  are *not* similar in any of the 20 bands:  $(1-0.328)^{20} = .00035$  .
  - i.e., about 1/3000th of the 80%-similar underlying sets are false negatives.

# LSH Summary

- Tune to get almost all pairs with similar signatures, but eliminate most pairs that do not have similar signatures.
- Check that candidate pairs really do have similar signatures.
- **Optional**: In another pass through data, check that the remaining candidate pairs really represent similar *sets* .

# Applications of LSH

**Entity Resolution**  
**Fingerprints**  
**Similar News Articles**

**Mining of Massive Datasets**  
**Leskovec, Rajaraman, and Ullman**  
**Stanford University**



# Entity Resolution

- The *entity-resolution* problem is to examine a collection of records and determine which refer to the same entity.
  - *Entities* could be people, events, etc.
- Typically, we want to merge records if their values in corresponding fields are similar.

# Matching Customer Records

- I once took a consulting job solving the following problem:
  - Company A agreed to solicit customers for Company B, for a fee.
  - They then argued over how many customers.
  - Neither recorded exactly which customers were involved.



# Customer Records – (2)

- Each company had about 1 million records describing customers that might have been sent from A to B.
- Records had name, address, and phone, but for various reasons, they could be different for the same person.

# Customer Records – (3)

- **Step 1:** Design a measure (“*score*”) of how similar records are:
  - E.g., deduct points for small misspellings (“Jeffrey” vs. “Jeffery”) or same phone with different area code.
- **Step 2:** Score all pairs of records that the LSH scheme identified as candidates; report high scores as matches.

# Customer Records – (4)

- **Problem:**  $(1 \text{ million})^2$  is too many pairs of records to score.
- **Solution:** A simple LSH.
  - Three hash functions: exact values of name, address, phone.
    - Compare iff records are identical in at least one.
  - Misses similar records with a small differences in all three fields.

## Aside: Hashing Names, Etc.

- How do we hash strings such as names so there is one bucket for each string?
- **Answer:** Sort the strings instead.
- Another option was to use a few million buckets, and deal with buckets that contain several different strings.

## Aside: Validation of Results

- We were able to tell what values of the scoring function were reliable in an interesting way.
- Identical records had a creation date difference of 10 days.
- We only looked for records created within 90 days of each other, so bogus matches had a 45-day average.

# Validation – (2)

- By looking at the pool of matches with a fixed score, we could compute the average time-difference, say  $x$ , and deduce that fraction  $(45-x)/35$  of them were valid matches.
- Alas, the lawyers didn't think the jury would understand.

# Validation – Generalized

- Any field not used in the LSH could have been used to validate, provided corresponding values were closer for true matches than false.
- **Example:** if records had a **height** field, we would expect true matches to be close, false matches to have the average difference for random people.

# Backup



# Implementation of Minhashing

- Suppose 1 billion rows.
- Hard to pick a random permutation of 1...billion.
- Representing a random permutation requires 1 billion entries.
- Accessing rows in permuted order leads to thrashing.

# Implementation – (2)

- A good approximation to permuting rows: pick, say, 100 hash functions.
- For each column  $c$  and each hash function  $h_i$ , keep a “slot”  $M(i, c)$ .
- **Intent:**  $M(i, c)$  will become the smallest value of  $h_i(r)$  for which column  $c$  has 1 in row  $r$ .
  - I.e.,  $h_i(r)$  gives order of rows for  $i^{\text{th}}$  permutation.

# Implementation – (3)

```
for each row  $r$  do begin  
  for each hash function  $h_i$  do  
    compute  $h_i(r)$ ;  
  for each column  $c$   
    if  $c$  has 1 in row  $r$   
      for each hash function  $h_i$  do  
        if  $h_i(r)$  is smaller than  $M(i, c)$  then  
           $M(i, c) := h_i(r)$ ;  
end;
```

# Example

Row	C1	C2
1	1	0
2	0	1
3	1	1
4	1	0
5	0	1

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

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

$$h(1) = 1 \quad \text{Sig1: } 1 \quad \text{Sig2: } \infty$$

$$g(1) = 3 \quad \text{Sig1: } 3 \quad \text{Sig2: } \infty$$

$$h(2) = 2 \quad 1 \quad \text{Sig2: } 2$$

$$g(2) = 0 \quad 3 \quad \text{Sig2: } 0$$

$$h(3) = 3 \quad 1 \quad 2$$

$$g(3) = 2 \quad \text{Sig1: } 2 \quad 0$$

$$h(4) = 4 \quad 1 \quad 2$$

$$g(4) = 4 \quad 2 \quad 0$$

$$h(5) = 0 \quad 1 \quad \text{Sig2: } 0$$

$$g(5) = 1 \quad 2 \quad 0$$

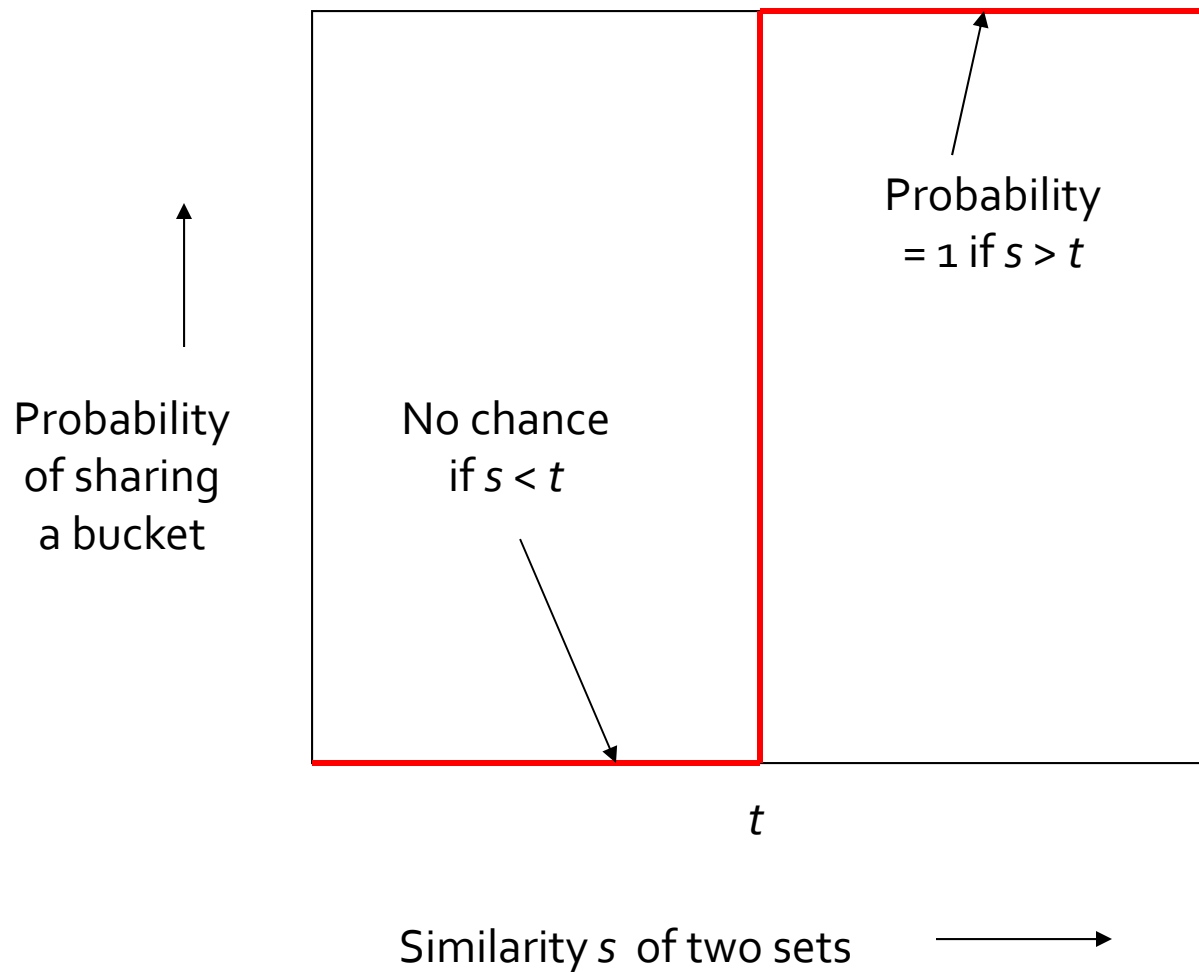
# Implementation – (4)

- Often, data is given by column, not row.
  - **Example**: columns = documents, rows = shingles.
- If so, sort matrix once so it is by row.

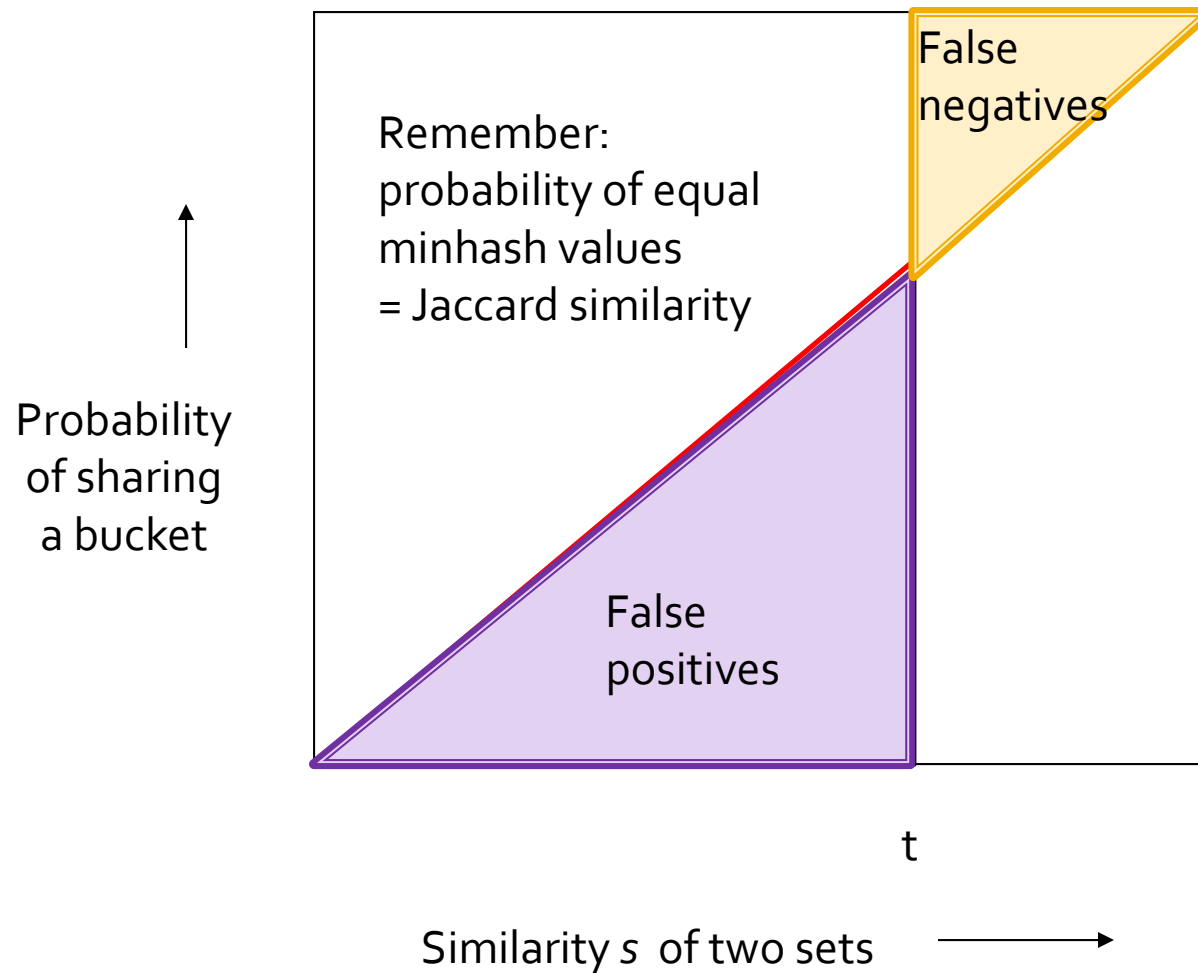
# Suppose $C_1, C_2$ Only 40% Similar

- Probability  $C_1, C_2$  identical in any one particular band:  $(0.4)^5 = 0.01$  .
- Probability  $C_1, C_2$  identical in  $\geq 1$  of 20 bands:  $\leq 20 * 0.01 = 0.2$  .
- But false positives much lower for similarities  $\ll 40\%$ .

# Analysis of LSH – What We Want

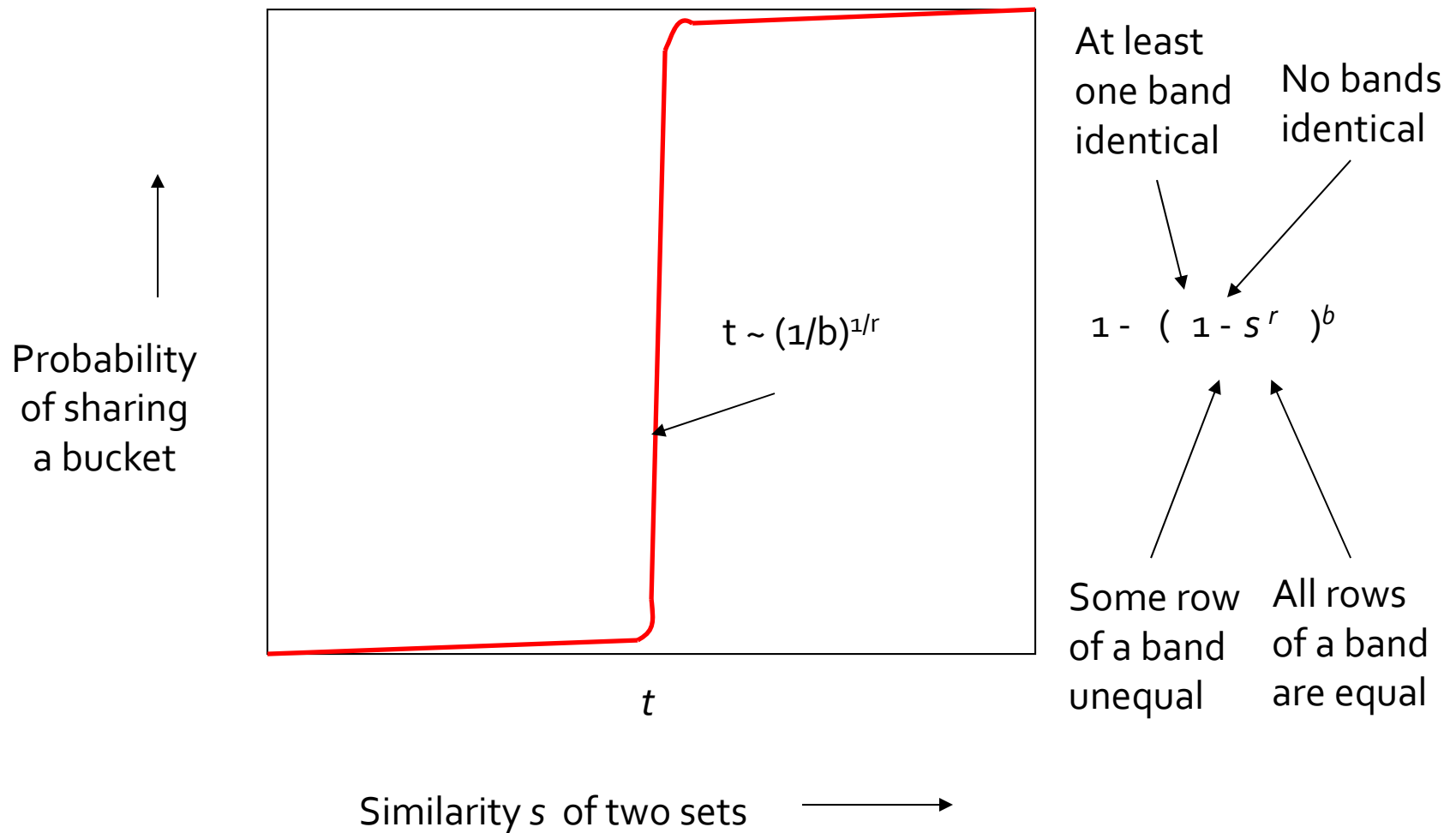


# What One Band of One Row Gives You





# What $b$ Bands of $r$ Rows Gives You



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

$s$	$1-(1-s^r)^b$
.2	.006
.3	.047
.4	.186
.5	.470
.6	.802
.7	.975
.8	.9996