

APPENDIX: HASH FUNCTIONS

(From Kleinberg & Tardos Book: Algorithm Design)

Dictionary data type

Dictionary. Given a universe U of possible elements, maintain a subset $S \subseteq U$ so that **inserting**, **deleting**, and **searching** in S is efficient.

Dictionary interface.

- $\text{create}()$: initialize a dictionary with $S = \emptyset$.
- $\text{insert}(u)$: add element $u \in U$ to S .
- $\text{delete}(u)$: delete u from S (if u is currently in S).
- $\text{lookup}(u)$: is u in S ?

Challenge. Universe U can be extremely large so defining an array of size $|U|$ is infeasible.

Applications. File systems, databases, Google, compilers, checksums, P2P networks, associative arrays, cryptography, web caching, etc.

HASH FUNCTIONS

Hashing

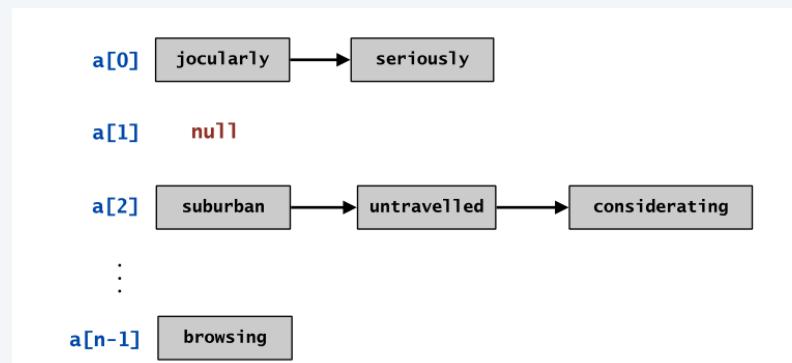
Hash function. $h : U \rightarrow \{ 0, 1, \dots, n - 1 \}$.

Hashing. Create an array a of length n . When processing element u , access array element $a[h(u)]$.

Collision. When $h(u) = h(v)$ but $u \neq v$.

birthday paradox
↙

- A collision is expected after $\Theta(\sqrt{n})$ random insertions.
- Separate chaining: $a[i]$ stores linked list of elements u with $h(u) = i$.



HASH FUNCTIONS

Hashing performance $h: U \rightarrow [n]; |U| \gg n; \forall S \subseteq U \text{ s.t. } |S| = m$

Ideal hash function. Maps m elements uniformly at random to n hash slots.

- Running time depends on length of chains.
- Average length of chain = $\alpha = m / n$.
- Choose $n \approx m \Rightarrow$ expect $O(1)$ per insert, lookup, or delete.

Challenge. Hash function h that achieves $O(1)$ per operation.

Approach. Use randomization for the choice of h .

adversary knows the randomized algorithm you're using,
but doesn't know random choice that the algorithm makes

Universal Family of Hash Functions

Universal hashing (Carter-Wegman 1980s)

A **universal family of hash functions** is a set of hash functions H mapping a universe U to the set $\{0, 1, \dots, n-1\}$ such that

- For any pair of elements $u \neq v$: $\Pr_{h \in H} [h(u) = h(v)] \leq 1/n$
- Can select random h efficiently.
- Can compute $h(u)$ efficiently.

chosen uniformly at random

Ex. $U = \{a, b, c, d, e, f\}$, $n = 2$.

	a	b	c	d	e	f
$h_1(x)$	0	1	0	1	0	1
$h_2(x)$	0	0	0	1	1	1

$$H = \{h_1, h_2\}$$

$$\Pr_{h \in H} [h(a) = h(b)] = 1/2$$

not universal

$$\Pr_{h \in H} [h(a) = h(c)] = 1$$

$$\Pr_{h \in H} [h(a) = h(d)] = 0$$

...

$$H = \{h_1, h_2, h_3, h_4\}$$

$$\Pr_{h \in H} [h(a) = h(b)] = 1/2$$

universal

$$\Pr_{h \in H} [h(a) = h(c)] = 1/2$$

$$\Pr_{h \in H} [h(a) = h(d)] = 1/2$$

$$\Pr_{h \in H} [h(a) = h(e)] = 1/2$$

$$\Pr_{h \in H} [h(a) = h(f)] = 0$$

...

	a	b	c	d	e	f
$h_1(x)$	0	1	0	1	0	1
$h_2(x)$	0	0	0	1	1	1
$h_3(x)$	0	0	1	0	1	1
$h_4(x)$	1	0	0	1	1	0

Universal Family of Hash Functions

Universal hashing: analysis

Proposition. Let H be a universal family of hash functions mapping a universe U to the set $\{0, 1, \dots, n - 1\}$; let $h \in H$ be chosen uniformly at random from H ; let $S \subseteq U$ be a subset of size at most n ; and let $u \notin S$. Then, the expected number of items in S that collide with u is at most 1.

Pf. For any $s \in S$, define random variable $X_s = 1$ if $h(s) = h(u)$, and 0 otherwise. Let X be a random variable counting the total number of collisions with u .

$$E_{h \in H}[X] = E[\sum_{s \in S} X_s] = \sum_{s \in S} E[X_s] = \sum_{s \in S} \Pr[X_s = 1] \leq \sum_{s \in S} \frac{1}{n} = |S| \frac{1}{n} \leq 1$$

↑
linearity of expectation ↑
 X_s is a 0–1 random variable ↑
universal

Q. OK, but how do we design a universal class of hash functions?

Universal Family of Hash Functions

Designing a universal family of hash functions

Modulus. We will use a prime number p for the size of the hash table.

Integer encoding. Uniquely identify each element $u \in U$ with a base- p integer of r digits: $x = (x_1, x_2, \dots, x_r)$.

distinct elements have
different encodings

Hash function. Let $A =$ set of all r -digit, base- p integers. For each $a = (a_1, a_2, \dots, a_r)$ where $0 \leq a_i < p$, define

$$h_a(x) = \left(\sum_{i=1}^r a_i x_i \right) \bmod p \quad \leftarrow \text{maps universe } U \text{ to set } \{0, 1, \dots, p-1\}$$

Hash function family. $H = \{ h_a : a \in A \}$.

Universal Family of Hash Functions

$$U = [p]^r$$

Designing a universal family of hash functions

$$h_a(x) = \sum_{i=1}^r a_i x_i \pmod{p}$$

Theorem. $H = \{ h_a : a \in A \}$ is a universal family of hash functions.

Pf. Let $x = (x_1, x_2, \dots, x_r)$ and $y = (y_1, y_2, \dots, y_r)$ encode two distinct elements of U .

We need to show that $\Pr[h_a(x) = h_a(y)] \leq 1/p$.

- Since $x \neq y$, there exists an integer j such that $x_j \neq y_j$.
- We have $h_a(x) = h_a(y)$ iff

$$(1) \quad a_j \underbrace{(y_j - x_j)}_z \equiv \underbrace{\sum_{i \neq j} a_i (x_i - y_i)}_m \pmod{p}$$

Principle of Deferred Decision
This is fixed before
 a_j is chosen i.u.a.r.

- Can assume a was chosen uniformly at random by first selecting all coordinates a_i where $i \neq j$, then selecting a_j at random. Thus, we can assume a_i is fixed for all coordinates $i \neq j$.
- Since p is prime, $a_j z \equiv m \pmod{p}$ has at most one solution among p possibilities. ← see lemma on next slide
- Thus $\Pr[h_a(x) = h_a(y)] \leq 1/p$. \rightarrow it is the a_j value s.t. (1) is satisfied

Universal Family of Hash Functions

Number theory fact

Fact. Let p be prime, and let $z \not\equiv 0 \pmod{p}$. Then $\alpha z \equiv m \pmod{p}$ has at most one solution $0 \leq \alpha < p$.

Pf.

- Suppose $0 \leq \alpha_1 < p$ and $0 \leq \alpha_2 < p$ are two different solutions.
- Then $(\alpha_1 - \alpha_2)z \equiv 0 \pmod{p}$; hence $(\alpha_1 - \alpha_2)z$ is divisible by p .
- Since $z \not\equiv 0 \pmod{p}$, we know that z is not divisible by p .
- It follows that $(\alpha_1 - \alpha_2)$ is divisible by p .
- This implies $\alpha_1 = \alpha_2$. ■

here's where we
use that p is prime

Bonus fact. Can replace “at most one” with “exactly one” in above fact.

Pf idea. Euclid’s algorithm.

Universal Family of Hash Functions

Universal hashing: summary

Goal. Given a universe U , maintain a subset $S \subseteq U$ so that insert, delete, and lookup are efficient.

Universal hash function family. $H = \{ h_a : a \in A \}$.

$$h_a(x) = \left(\sum_{i=1}^r a_i x_i \right) \bmod p$$

- Choose p prime so that $m \leq p \leq 2m$, where $m = |S|$.
- Fact: there exists a prime between m and $2m$. ← can find such a prime using another randomized algorithm (!)

Consequence.

- Space used = $\Theta(m)$.
- Expected number of collisions per operation is ≤ 1
 $\Rightarrow O(1)$ time per insert, delete, or lookup.

Note to other teachers and users of these slides: We would be delighted if you found this our material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. If you make use of a significant portion of these slides in your own lecture, please include this message, or a link to our web site: <http://www.mmds.org>

Finding Similar Items in Large Data Sets

Lecturer: **Andrea Clementi**
University of Rome *Tor Vergata*

Note: These slides are an adaptation, with more details, of the slides by Jure Leskovec, Anand Rajaraman, Jeff Ullman available at <http://www.mmds.org>

New thread: High dim. data

High dim.
data

Locality
sensitive
hashing

Clustering

Dimensional
ity
reduction

Graph
data

PageRank,
SimRank

Network
Analysis

Spam
Detection

Infinite
data

Filtering
data
streams

Web
advertising

Queries on
streams

Machine
learning

SVM

Decision
Trees

Perceptron,
kNN

Apps

Recommen
der systems

Association
Rules

Duplicate
document
detection

Application: Scene Completion Problem

target
item



many
items



selection of
similar items

closest
item



A Common Metaphor

Many problems can be expressed as:

finding “similar” sets

Examples:

- **Pages with a large fraction of similar words**
 - For duplicate detection, classification by topic
- **Customers who purchased a large fraction of similar products**
 - Products with similar customer sets
- **Images with similar features**
 - Users who visited similar websites

High-Dimensional Data and Distance

Input: *High-Dimensional* data points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$

- Image \mathbf{x} is a *long* vector of pixel colors

$$\mathbf{x} = \begin{bmatrix} 1 & 4 & 1 \\ 0 & 2 & 1 \\ 0 & 1 & 0 \end{bmatrix} \rightarrow [1 \ 4 \ 1 \ 0 \ 2 \ 1 \ 0 \ 1 \ 0] \in \{0,1,\dots,c\}^h, \quad (h = \text{data dimension})$$

Introduce a **distance function** $d(\mathbf{x}_1, \mathbf{x}_2)$

- Which quantifies the “distance” between any pair of data points \mathbf{x}_1 and \mathbf{x}_2

Goal: Find ***all pairs of data points*** $(\mathbf{x}_i, \mathbf{x}_j)$ that are within some *distance threshold* $d(\mathbf{x}_i, \mathbf{x}_j) \leq s$

High-Dimensional Data and Distance

Goal: Find **all pairs of data points** (x_i, x_j) that are within some distance threshold $d(x_i, x_j) \leq s$

Naïve solution would take $O(N^2 * h)$ 😞

where N is the number of data points and h is the data dimension.

MAGIC: This can be done in $O(N * h')$ with $h' \ll h$, **with high confidence !!**

Why? How?

Our Application: Finding Similar Documents

Goal: Given a large number (N in the millions or billions) of **Docs**, find “near duplicate” pairs

Applications:

- Mirror websites, or approximate mirrors
 - Don’t want to show both in search results
- Similar news articles at many news sites
 - Cluster articles by “same story”

Technical Problems:

- Many small *pieces* of one **Doc** can appear out of order in another
- Too many **Docs** to compare all pairs
- **Docs** are so large or so many that they cannot fit in main memory

Choosing the *right* Distance Measures

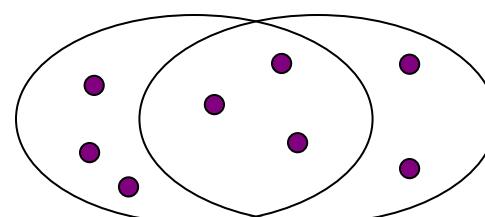
For each application, we first need to define what “**distance**” means

Our Choice: *Jaccard* distance/similarity
over a family of subsets

- The *Jaccard* similarity of two **sets** is:

$$sim(C_1, C_2) = |C_1 \cap C_2| / |C_1 \cup C_2|$$

- Jaccard distance:** $d(C_1, C_2) = 1 - sim(C_1, C_2)$



3 in intersection
8 in union
Jaccard similarity= 3/8
Jaccard distance = 5/8

Application: Documents

A First Important Question: How apply *Jaccard distance/similarity* to **Docs** (i.e. to finite-size *strings*)?

Doc $\in U^*$ where **U** is the **Alphabet**

$$J.Similarity \quad sim(C_1, C_2) = |C_1 \cap C_2| / |C_1 \cup C_2|$$

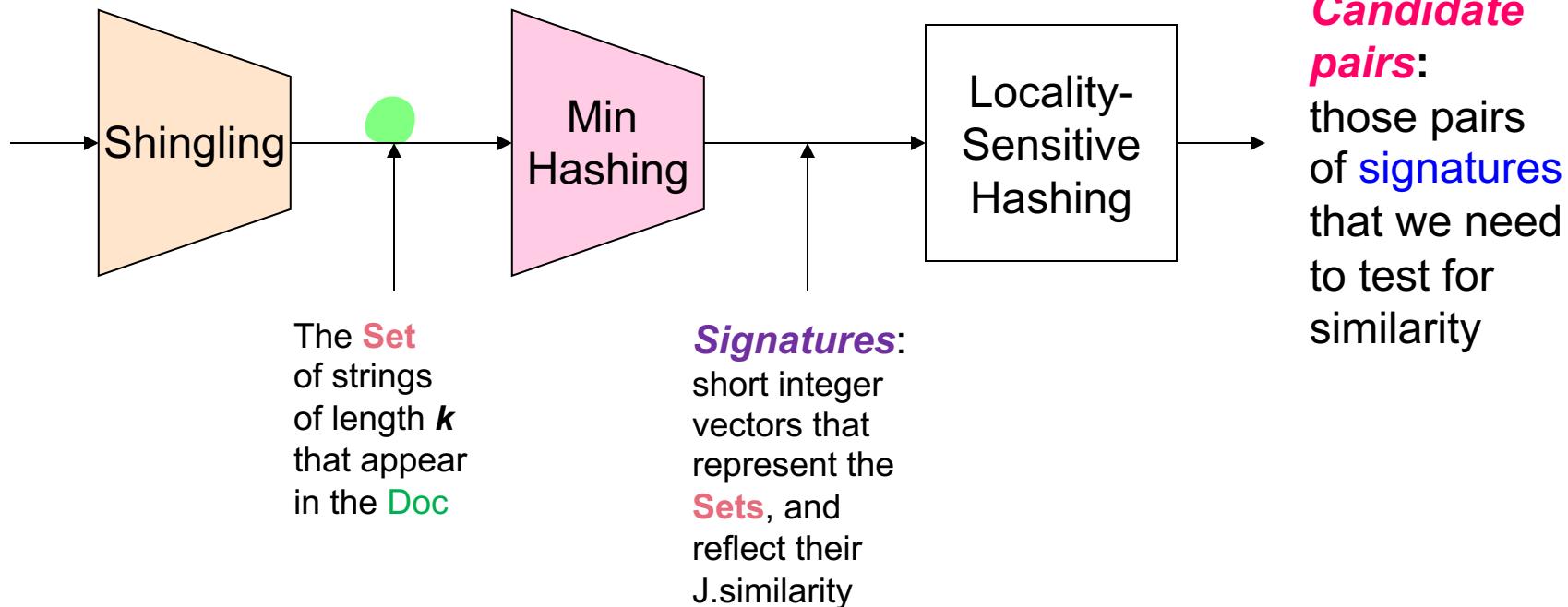
Docs \Leftrightarrow **Sets** ? 

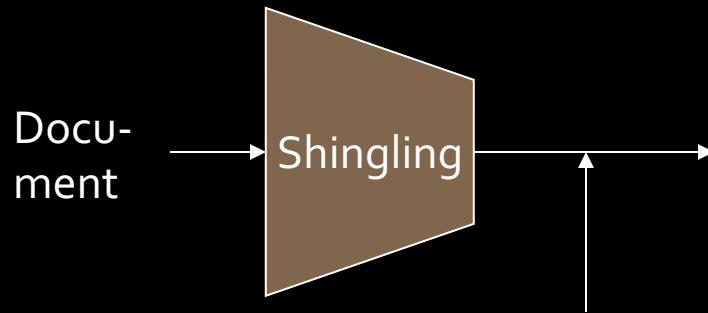
3 Essential Steps for Similar Docs

1. ***Input:*** A Huge Universe of *Docs*
2. ***Shingling:*** Convert *Docs* to (large) *Sets*
3. ***Min-Hashing:*** Convert large *Sets* to short *Signatures*, while preserving ***J. Similarity***
4. ***Locality-Sensitive Hashing:*** Focus on pairs of *Signatures* likely to be from *J.similar Docs* (overcome the $\Theta(N^2)$ barrier!!)
5. ***Output: Candidate Doc pairs!***

The Big Picture

Doc





The set
of strings
of length k
that appear
in the doc-
ument

Shingling

Step 1: *Shingling*: Convert documents to sets

Documents as High-Dim Data

Step 1: *Shingling*: Convert Documents to Sets

Simple approaches:

- Document = set of *words* appearing in document
- Document = set of “important” *words*

Don’t work well for this application.....



Why?

Need to account for ordering of words!

A different approach: Shingles!

Definition of Shingles

A *k-shingle* (or *k-gram*) for a Doc is a sequence of *k tokens* that appears in the Doc

- Tokens can be *characters*, *words* or something else, depending on the application:
 - $U \equiv$ Universe of *all* possible Tokens
- Example: Assume Tokens = **characters** and **k=2**
 - Then for Doc $D_1 = abcab \rightarrow$ Set of *2-shingles* is $S(D_1) = \{ab, bc, ca\}$
- Option: Shingles as a bag (multiset). Count **ab twice**:
 $S'(D_1) = \{ab, bc, ca, ab\}$

Doc Representation by Shingles

Each Doc $D \in U^*$ is represented as the set of its **k-shingles** $C=S(D)$

Equivalently:

each Doc is represented as a *long binary vector* $C=S(D)$ having one component for each element of the set U^k of all possible **k-shingles**

- Each unique shingle is a *dimension (component)* of $C=S(D)$
- **Note:** Vectors have several dimensions but are very sparse

Example:



Similarity Metric for Shingles

Our strategy:

(a) Doc D is represented as the subset of its **k-shingles** $C_j = S(D)$

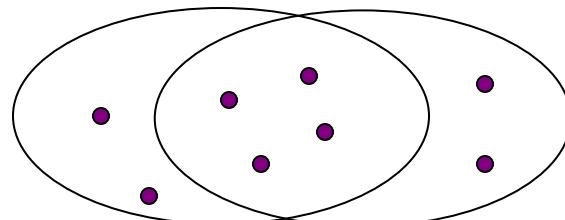
Equivalently, each Doc is a **binary vector** in the space of all possible k -shingles U^k

(b) A natural **similarity** measure (=1- distance) is the **Jaccard similarity** (and, thus, **distance**):

$$J.sim(D_1, D_2) \equiv |C_1 \cap C_2| / |C_1 \cup C_2| = \Pr(C_1 \cap C_2 | C_1 \cup C_2)$$

PROBABILISTIC
INTERPRETATION

$$\text{dist}(D_1, D_2) \equiv 1 - sim(D_1, D_2)$$



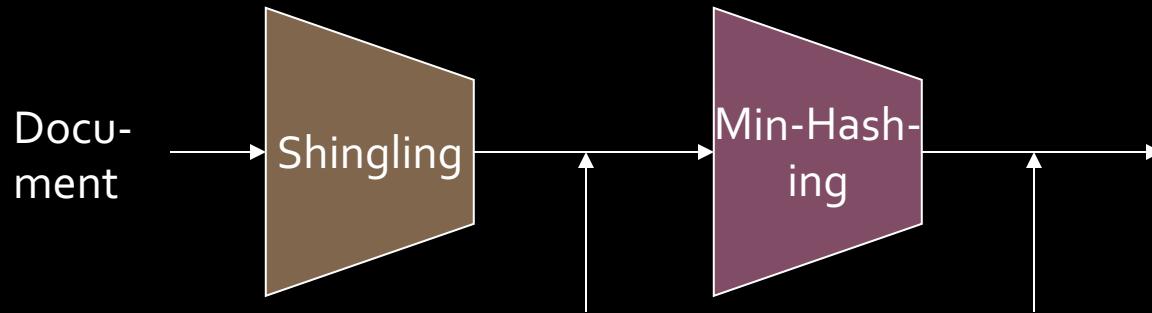
Working Hypothesis

Docs that have lots of k-shingles in common
have similar text, even if the text appears in
different order

Practice: You must pick k large enough, or most Docs
will have most shingles → large intersection →
large (not relevant) *J. Sim.*

Parameter tuning from Real World:

- $k = 5$ is OK for short Docs
- $k = 10$ is better for long Docs



The set
of strings
of length k
that appear
in the doc-
ument

Signatures:
short integer
vectors that
represent the
sets, and
reflect their
similarity

MinHashing

Step 2: *Min-hashing*: Convert **large sets** (i.e. long binary sparse vectors) to short **signatures**, while **preserving similarity**

Recall: Motivation for Min-Hashing

Suppose we need to find near-duplicate **Docs** among $N = \cancel{10^6} \cdot 10^6$ **Docs** and, hence, among $\cancel{10^6} \cdot 10^6$ **k-shingles** $C = S(D)$

Then....

$$\approx \frac{1}{2} \cdot 10^{6+6}$$

Min-Hashing does not avoid $\Theta(N^2) \simeq 5 \cdot 10^{11}$ pair comparisons!

However ...

It makes each pair comparison (*J. similarity* computations) much more efficient



Note: if $C \in \{0,1\}^{U^k} \rightarrow \text{cost(Naïve-Comparison)} = \Theta(|U|^k)$. For alphabet size $|U| = 27$ and $k = 10 \rightarrow \dots$ 😱

Encoding Sets as Binary Vectors

Many **similarity** problems can be formalized as **finding subsets** that have significant intersection (recall *J. Sim.*):

(a) Encode sets using binary vectors

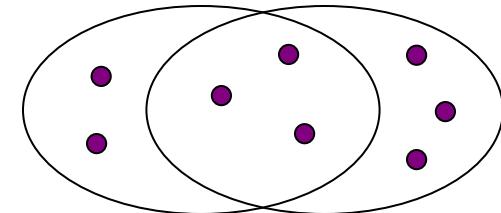
- One dimension (component) per *element* (**shingle**) in the *Universe* U^k

(b) Interpret:

- (i) *set intersection* as bitwise AND
- (ii) *set union* as bitwise OR

Example: $C_1 = 10111$; $C_2 = 10011$

- Size of **intersection** = 3; size of **union** = 4,
- **Jaccard similarity** (not distance) = $3/4$
- **Distance:** $d(C_1, C_2) = 1 - (J.Sim(C_1, C_2)) = 1/4$



From Sets to Large, Sparse Boolean Matrices

Rows = elements (**shingles**) = vector

components (**n.** of rows = $m = |U|^k$)

Columns = Sets = Docs (**N col**)

- **1** in row s and column e iff $s \in e$
- Column similarity is the **J.sim** of the corresponding sets
- **Remark:** Typical matrix **M** is sparse!

Each Doc is a column:

- **Example:** $\text{sim}(C_1, C_2) = ?$

- Size of intersection = 3; size of union = 6,
Jaccard similarity (not distance) = $1/2$
- $d(C_1, C_2) = 1 - (\text{J.Sim } (C_1, C_2)) = 1/2$

Matrix **M**:

$|Docs| \geq N$

1	1	1	0
1	1	0	1
0	1	0	1
0	0	0	1
1	0	0	1
1	1	1	0
1	0	1	0

$$m \equiv |U|^k$$



Outline: Finding Similar Columns

So far:

- Docs → *Subsets* of Shingles
- Represent *Subsets* as *binary columns* of a matrix \mathbf{M}

Next goal: Find similar columns by computing
small column signatures (**Min-Hashing**)

Technical goal: define a *good* signature *algorithm* so that
J.Sim of Columns ≈ “Similarity” of Signatures

Outline: Finding Similar Columns

Next Goal: Find similar **columns** using small signatures

Our Approach:

- 1) Signatures of **columns** = small *summaries* of **columns**
- 2) Examine pairs of **signatures** to find similar **columns**
 - Essential: Similarities of **signatures** and **columns** must be related
- 3) Optional: Check that **columns** with similar **signatures** are really *similar*

Warnings:

- Comparing *all pairs* may take too much time: **Task for LSH**
- {These methods can produce false negatives, and even false positives (if the optional check Step (3) is not made)}

not now

Hashing Columns (Signatures)

Key idea: “hash” each m -size column C to a small *signature* $h(C)$, such that: $|h(C)| \ll |C| \leq m$

- (1) $h(C)$ is small enough that the signature fits in RAM (Fast Memory)
- (2) $\text{sim}(C_1, C_2)$ is “close” to the *equivalence* of signatures $h(C_1)$ and $h(C_2)$:

Goal: Find a *hash* function $h(\cdot)$ such that:

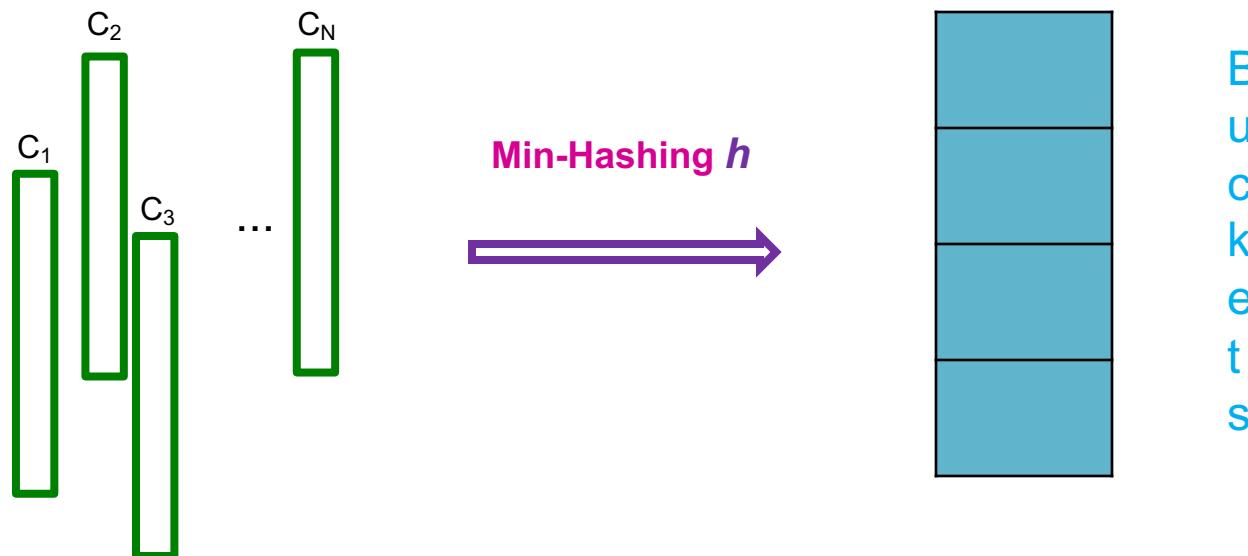
- If $\text{sim}(C_1, C_2)$ is high, then **with high prob.** $\underline{h(C_1) = h(C_2)}$
- If $\text{sim}(C_1, C_2)$ is low, then **with high prob.** $\underline{h(C_1) \neq h(C_2)}$

Hashing Columns (Signatures)

Goal: Find a hash function $h(\cdot)$ such that:

- If $\text{sim}(C_1, C_2)$ is high, then with high prob. $h(C_1) = h(C_2)$
- If $\text{sim}(C_1, C_2)$ is low, then with high prob. $h(C_1) \neq h(C_2)$

Intuition: Hash Docs into *Buckets*. Expect that “most” pairs of near duplicate Docs hash into the same Bucket!



Min-Hashing

Goal: Find a hash function $h(\cdot)$ such that:

- If $\text{sim}(C_1, C_2)$ is high, then with high prob. $h(C_1) = h(C_2)$
- If $\text{sim}(C_1, C_2)$ is low, then with high prob. $h(C_1) \neq h(C_2)$

Clearly, the hash function depends on
the *similarity* metric:

- Not all similarity metrics have a suitable
hash function

Good News: There is a suitable hash function for
the Jaccard similarity: It is called Min-Hashing

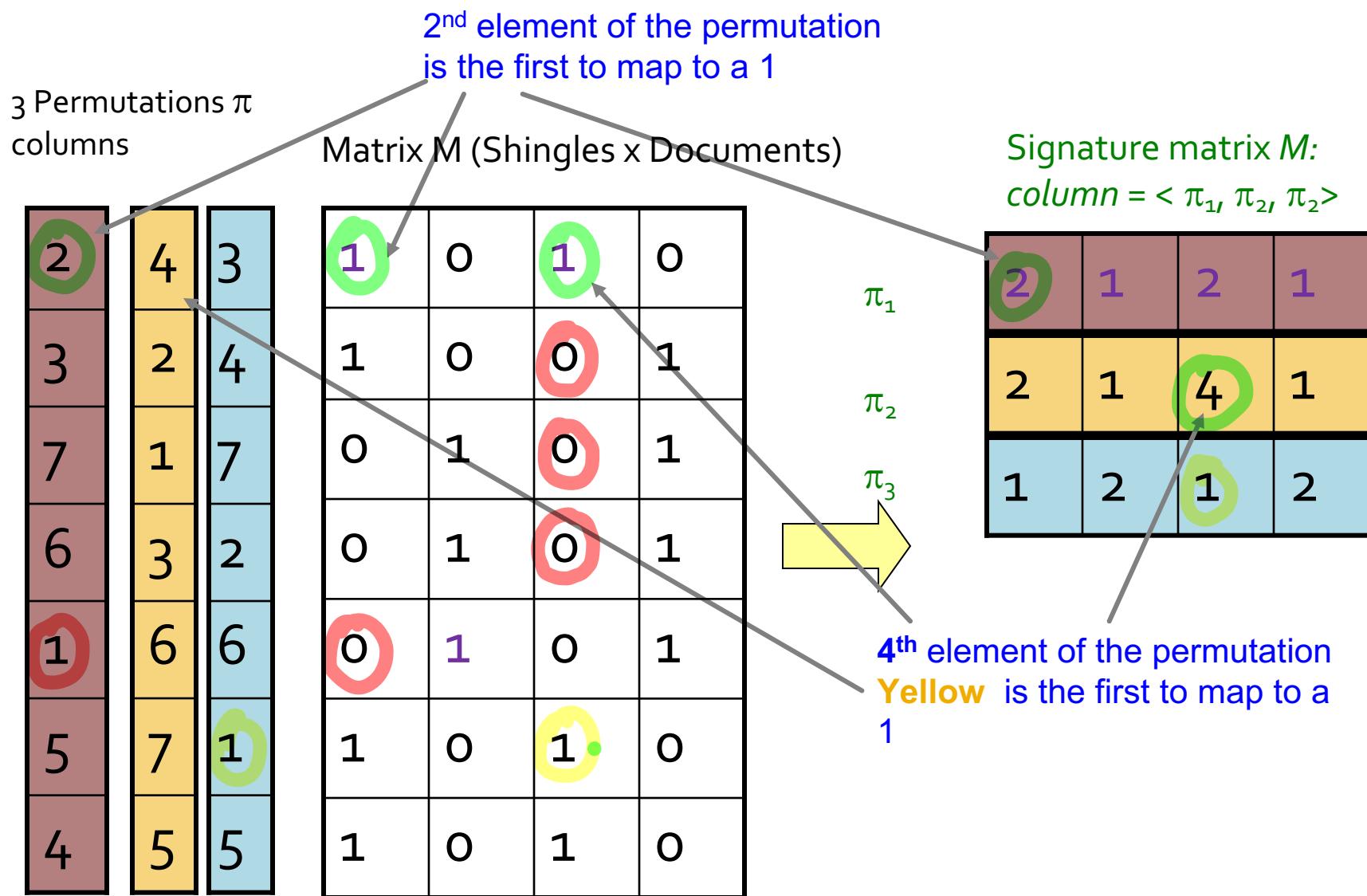
Min-Hashing

Imagine the rows of matrix \mathbf{M} permuted under *random permutation* π $\mathbf{M} =$ ORIGINAL MATRIX

Def. For each column C , define: **Min-Hash function** $h_{\pi}(C) \equiv$ the index of the first (in the permuted order π) row in which column C has value 1 (we named it as $\min(\pi(C))$)

Practical Strategy: Use several (e.g., 100) independent rnd **hash** functions (that is, rnd permutations) to create a *signature* of C
see Cates....

Min-Hashing Example with N=4, m= 7



MEMORY - COST (MIN-HASH)

$$h_{\pi}(c) \rightarrow [m] = \{1, 2, \dots, m\}$$

where $|C| \equiv |U|^k \equiv m$

and $\text{SPACE}(h_{\pi}(c)) \in \Theta(\lg m)$

$$m \xrightarrow{h_{\pi}} \Theta(\lg m)$$

m \downarrow
it is just a
row-index of M

The Min-Hash Property

3-SIMILARITY PRESERVING

THM: Fix C_1, C_2 . Choose u.a.r. a permutation π . Then,

$$\Pr_{\pi} [h_{\pi}(C_1) = h_{\pi}(C_2)] = J.sim(C_1, C_2)$$

$$\frac{1}{3}$$


Proof.

- Let X be a **Doc** (column of shingles), $y \in X$ is one of its **shingle**. Then: $\Pr[\pi(y) = \min(\pi(X))] = 1/|X|$ since:
 - It is equally likely that any row $y \in X$ is mapped to the **min** index
- Let y be s.t. $\pi(y) = \min(\pi(C_1 \cup C_2))$ π is rnd unif
- Then either: $\pi(y) = \min(\pi(C_1))$ if $y \in C_1$, or → UNION
 - $\pi(y) = \min(\pi(C_2))$ if $y \in C_2$ One of the two cols had to have
 - at position y
- So the prob. that **both** are **true** is the prob. of $y \in C_1 \cap C_2$
- $\Pr[\min(\pi(C_1)) = \min(\pi(C_2))] = |C_1 \cap C_2| / |C_1 \cup C_2| = sim(C_1, C_2)$

$\underbrace{h(C_1)}_{h(C_1)}$
 $\underbrace{h(C_2)}_{h(C_2)}$
 \hookrightarrow

(
| $C_1 \cap C_2| \cdot \frac{1}{|X|}$

0	0
0	0
1	1
0	0
0	1
1	0

45

Similarity for Signatures

Given any two columns C_1 and C_2 , Main THM \rightarrow the Prob. that C_1 and C_2 get the same Min-Hash signature, i.e. $h_\pi(C_1) = h_\pi(C_2)$, equals the value of $J.\text{sim}(C_1, C_2)$. Increase the "Confidence" ?? How ?

Simple Idea from Statistics: Exploit concentration on the Expected Value: use multiple, mutually-independent Min-Hash signatures $SIG(C) = \langle h_{\pi_1}(C), h_{\pi_2}(C), \dots, h_{\pi_t}(C) \rangle$ for some fixed $t >> 1$ and set:

DEF. $\text{Sign-Sim}(C_1, C_2)$ of two (vector) signatures $SIG(C_1)$ and $SIG(C_2)$ = the fraction of the chosen Min-Hash signatures in which they agree

AVERAGE $\xrightarrow{t \rightarrow +\infty}$ EXPECTATION

Formally,

Def. binary r.v. $Z = 1$ iff " $\min(\pi(C_1)) = \min(\pi(C_2))$ ".

Goal. Estimate $\Pr[Z=1]$

$$\Pr[Z=1] \approx S.\text{sim}(C_1, C_2)$$

Min-Hashing Example

\mathbb{M}

Permutation π

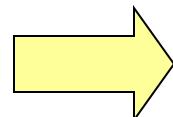
2	4	3
3	2	4
7	1	7
6	3	2
1	6	6
5	7	1
4	5	5

Input matrix (Shingles x Documents)

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 $SIG(\pi, C)$

6 rows
11
3



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

↳ Similarities: $|C_i \cap C_j| / |C_i \cup C_j|$

Col/Col
Sig/Sig

1-3	2-4	1-2	3-4
6/7	0.75	0	0
2/3	1.00	0	0

Min-Hash: The Algorithm

Rnd Algorithm *Doc-Pair Check*

Input: Columns $C_1, C_2 \in \{0,1\}^h$; Confidence Parameter t

For $j=1$ to t

- Choose u.a.r. a raw permutation $\pi_j \in \Pi_m$
- Compute hash $h_{\pi_j}(C_1)$ and $h_{\pi_j}(C_2)$

Return $\text{Sign-Sim}(C_1, C_2) = |\{j : h_{\pi_j}(C_1) = h_{\pi_j}(C_2)\}| / t$

Note: $J.Sim(C_1, C_2) \neq \text{Sign-Sim}(C_1, C_2)$ but... Because of the Min-Hash property (Main THM): $J.Sim(C_1, C_2) \equiv E_\pi(\text{Sign-Sim}(C_1, C_2))$, so

Coroll. 1 (of Main THM). $\text{Sign-Sim}(C_1, C_2) \rightarrow J.Sim(C_1, C_2)$ as $t \rightarrow \infty$

Important Note: since the t hash functions are mut. indep. \rightarrow we have concentration results on $|\text{Sign-Sim}(C_1, C_2) - J.Sim(C_1, C_2)|$

Min-Hash: Space Complexity

Pick $t=100$ random permutations of the m rows

Think of $SIG(C)$ as a short column vector of a signature Matrix

- $SIG(i,C)$ = the index of the first row that has a 1 in column C according to the i -th rnd permutation,

$SIG(i,C)$ is an index! $\Theta(\log |C|) = \Theta(\log m)$

Note: The sketch (signature) $SIG(*,C)$ of C is small: ~ 100 bytes
(in general: $\Theta(t * \log m)$)!

Recall: $m = |U|^k$ is the number of all possible **k-shingles**

We achieved our goal: We “compressed” long-bit vectors C into short signatures $SIG(*,C)$ ($m \rightarrow t \log m$)

Implementation Trick for Min-Hash Signatures

Min-Hash function $h_{\pi}(C)$ requires to generate several independent rnd permutations π of the raw set U^k ($|U|^k = m$)

Permuting m rows even once is too expensive ! 😞

Idea: Replace t Row Permutations with t Row Hash Functions

Ordering under an *Hash Function* f gives an (*almost*) random π

Pick $t = 100$ hash functions $f_i : [m] \rightarrow [m]$ (*there might be few collisions: neglect them!*)

Original Binary Matrix M

1	0	1	1	0	0
0	1	1	1	1	0
0	1	0	0	0	0
0	0	0	1	0	0
1	0	0	0	1	1
0	0	0	1	0	1

m rows = k-shingles

Min-Hash using f_i

Signature Matrix $SIG(i, C)$

f_1					
f_2					
...					
f_t					

N docs = columns

N docs = columns

raw index of the first 1
"after" hashing f_i

$\Theta(\lg m)$ bits

Algorithm for Signature Matrix $SIG(i, C)$

FOR each column $C=1\dots N$ and (hash-function) $i=1\dots t$ DO

- Initialize the signature matrix $SIG(i, C) = \infty$

FOR each row $j=1\dots m$ of the original matrix M DO

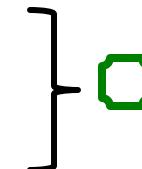
- FOR each column $C=1\dots N$ DO

- IF $M(j, C) = 1$ THEN (**)

- FOR each (hash-function) $i=1\dots t$ DO

- Compute the i -hash of raw j : $f_i(j)$

- IF $f_i(j) < SIG(i, C)$ THEN update $SIG(i, C) \leftarrow f_i(j)$



H(N - Computation
SEE PREVIOUS
LESSONS !

Performance Analysis: $\Theta(mN)$ iterations



Each iteration \square requires t hash computations $f_i(j)$ iff Condition $(**)$ holds.

Good News: i) No full computations of rnd permutations π_m ;

ii) Original Matrix M is sparse \rightarrow $(**)$ is unlike !

Bad News: Hash may generate **collisions** (no permutation)

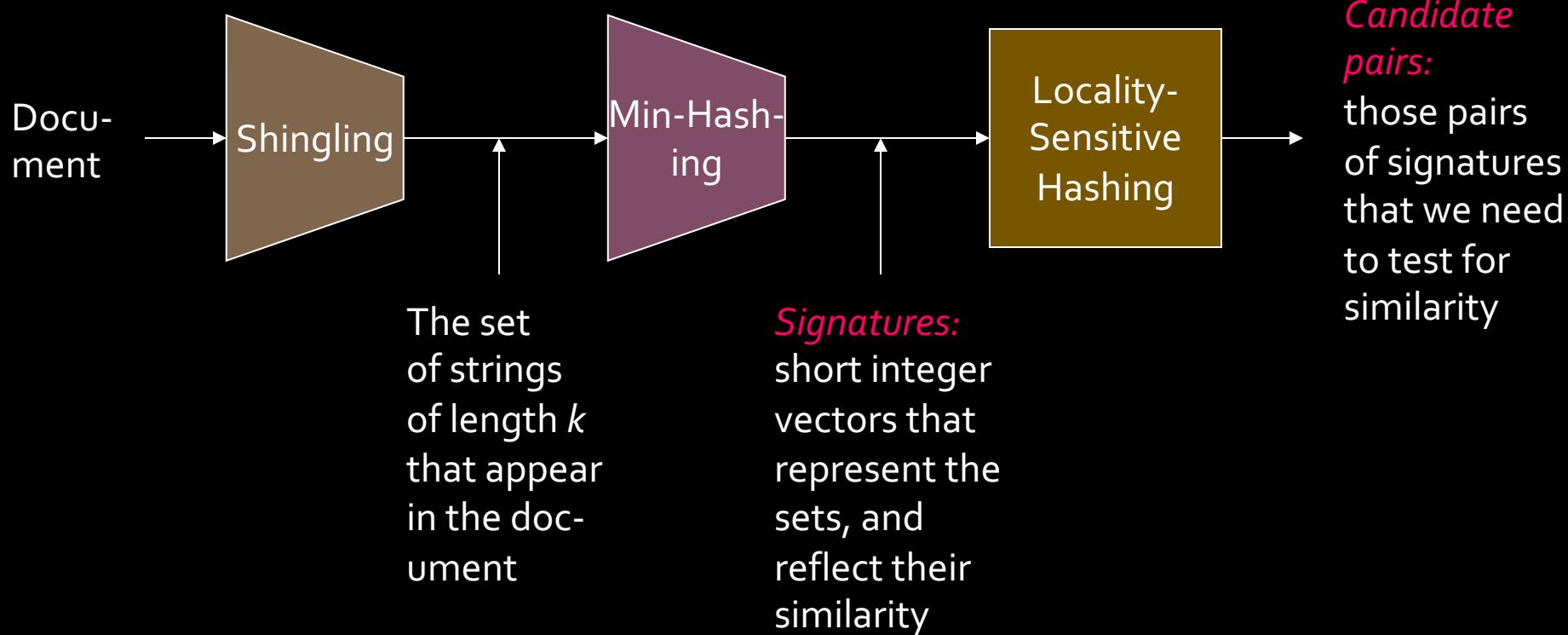
How to pick a random hash function $h(x)$?

Universal hashing:

$h_{a,b}(x) = ((a \cdot x + b) \bmod p) \bmod n$
where:

a, b ... random integers

p ... prime number ($p > m$)



Locality Sensitive Hashing

Step 3: *Locality-Sensitive Hashing:*

Focus on pairs of signatures likely to be from similar documents

Motivation for LSH

Suppose we need to find *near-duplicates* among $N = \cancel{10^6}$ Docs
 $\cancel{10^6}$

Naïvely, we would have to compute pairwise

J. similarities for every pair of Docs

$N(N - 1)/2 \approx 5 * 10^{11}$ comparisons !!!

- At 10^5 secs/day and 10^6 comparisons/sec, it would take **5 days**
- For $N = \cancel{10^6}$, it takes more than a year...!
 $\cancel{10^7}$

Remark: Min-Hashing does not avoid $\Theta(N^2)$ pair comparisons. It makes each pair comparison (J. similarity computations) much more efficient:



Local Sensitive Hashing (LSH)

Goal: Find **Doc** pairs with **similarity** at least **s** (for some **similarity threshold s** , e.g., $s = 0.8$)

L_p INPUT
PARAMETER

LSH – General idea: Use a function $f(x,y)$ that tells whether **x** and **y** is a **candidate Doc Pair**: a pair of **Docs** whose **similarity** must be evaluated in a deeper way.

General Idea of LSH: Hashing Docs Signature

Define **Candidate Function $f(x,y)$:**

- Consider (Min-Hash) Signature Matrix $SIG(i,C)$ ($i=1\dots t$; $C=1\dots N$)

$SIG(i,C)$

*

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

NOTE: From now on, we "forget" matrix M , and we only work with matrix $SIG(i,C)$ *

- Hash columns (Docs) of $SIG(i,C)$ to many Buckets
- Each pair of Docs that hashes into the same Bucket is a candidate pair: check it!

LSH
STRATEGY

LSH: Step I (Candidates from Min-Hash Signatures)

Choose a *suitable similarity threshold s ($0 < s < 1$)*

CRITERION

DEF. Columns x and y of $SIG(*, *)$ are a **candidate pair** if their Min-Hash signatures $SIG(*, x)$ and $SIG(*, y)$ agree on at least fraction s of their rows:
 $| \{i \in [t] : SIG(i, x) = SIG(i, y)\} | / t \geq s$

Working Assumption: Docs x and y have the same **similarity** as their **Signatures** (thanks to **Corollary 1 of Main THM**)

Note: Similarity of two column **signatures** does not come from *small differences* in some/several of their row **values**: there must be several/most **identical values**!

↓
index gap $|i-j|$ has
NO meaning!

LSH: Step II (Hash Signatures to Buckets)

Main idea: Hash Columns of $SIG(*, *)$ to a set B of Buckets for several times

min-hash
docs with same index
will be hashed to the same BUCKET

Key Fact: (only) candidate Column sets are likely to hash to the same Bucket, with high probability. (Key Property of Hash Functions)

Hence:

Candidate pairs are those that hash to the same Bucket

↳ This is a 2nd level of hashing

Exercise: LSH overcomes the $\Theta(N^2)$ barrier: Why?
When?

in
LSH

LSH. Step II (Details: Partition $SIG(*, *)$ into Bands)

Technical Problem: Hashing the full signature of a Doc into a one Bucket may generate several false positive and false negative.

A Good Solution.

Partition each doc signature, i.e. each column of $SIG(*, *)$ into a set of Bands (i.e. subsets of *adjacent* raws):

Divide matrix $SIG(*, *)$ into b Bands, each of r rows

For each Band, hash its raw-portion of each column to a hash table with c Buckets

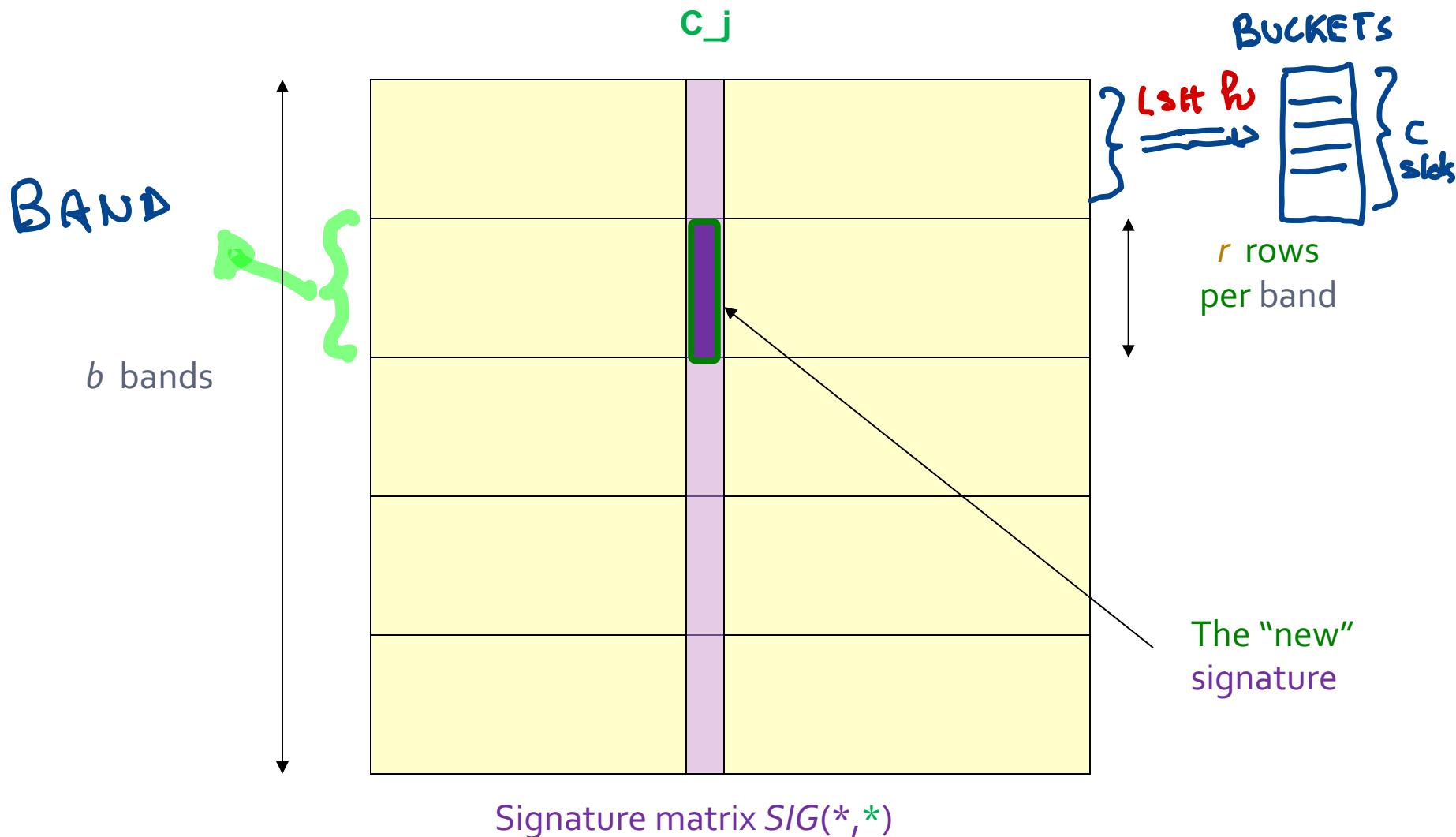
- Make c large enough (w.r.t. r) to *whp* avoid collisions of different signatures

[**Candidate Criterium:** Let's **candidate** column pairs those that hash to the same Bucket for at least one Band]

Practical Step: Tune b , r (*and thus t = total n. of raws*) to catch most similar pairs, but few non-similar pairs

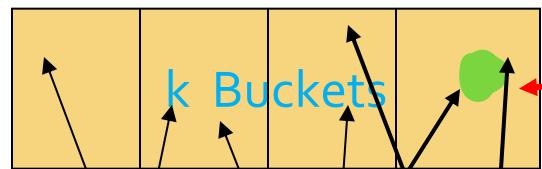
Partition $SIG(*, *)$ into b Bands

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



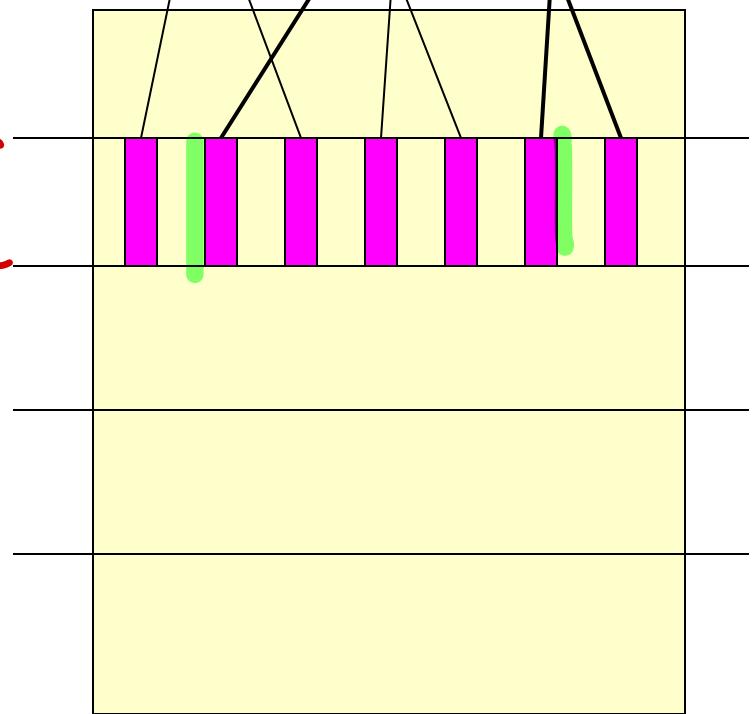
Hashing Bands

Each Band is
hashed to one
Bucket



Columns 2 and 6
are probably J. similar
(candidate pair)

BAND {



Columns 6 and 7 are
different.

r rows

b bands

LSH: Simplifying Assumptions

There are **enough Buckets** ($c \gg r$) that **Columns** are unlikely to hash to the same **Bucket** unless they are **identical** in a particular **Band**

We assume that: “**same Bucket**” means “**identical in that Band**” (Perfect Hashing)  (real app \approx universal hashing)

Assumptions above needed only to simplify analysis, not for correctness of algorithm

Recall: **candidate column** pairs = those that **hash** to the same Bucket for at least one **Band**

Example of Bands

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

Assume the following case:

Suppose 100,000 columns of $SIG(*, *)$ (100k Docs) $N = 10^5$

Signatures of 100 integers (rows) $\rightarrow t \leq 100$

Therefore, signatures take 40 Mb \equiv overall

Choose $b = 20$ bands, each of $r = 5$ integers.

Goal: Find pairs of Docs that are at least $s = 0.8$ similar

↳ INPUT parameter

Case I: C_1, C_2 are 80% similar

Find pairs of $\geq s=0.8$ similarity, set $b=20$, $r=5$

(i) Assume: $J.\text{sim}(C_1, C_2) = 0.8$

- Since $\text{sim}(C_1, C_2) \geq s$, we want C_1, C_2 to be a candidate pair: We want them to hash to at least one common Bucket (i.e. at least one Band is identical)

Fact 1. Probability C_1, C_2 get identical in one particular Band is

(i) $(0.8)^5 = 0.328$ {each of the 5 row pair must match} 

Fact 2. Probability C_1, C_2 get not identical in all of the 20 Bands is

(ii) $(1-0.328)^{20} = 0.00035$ {each of the 20 Bands does not satisfy (i)} false negative 

Confidence Rate: about 1/3000th of the 80%-similar column pairs are false negatives (we miss them) → We would find 99.965% pairs of truly similar Docs



Case II: C_1, C_2 are 30% similar

Find pairs of $\geq s=0.8$ similarity, set $b=20$, $r=5$

- (ii) Assume: $J.\text{sim}(C_1, C_2) = 0.3$
- Since $J.\text{sim}(C_1, C_2) \ll s$ we want C_1, C_2 to hash to **NO common Buckets** (*all Bands should be different*)

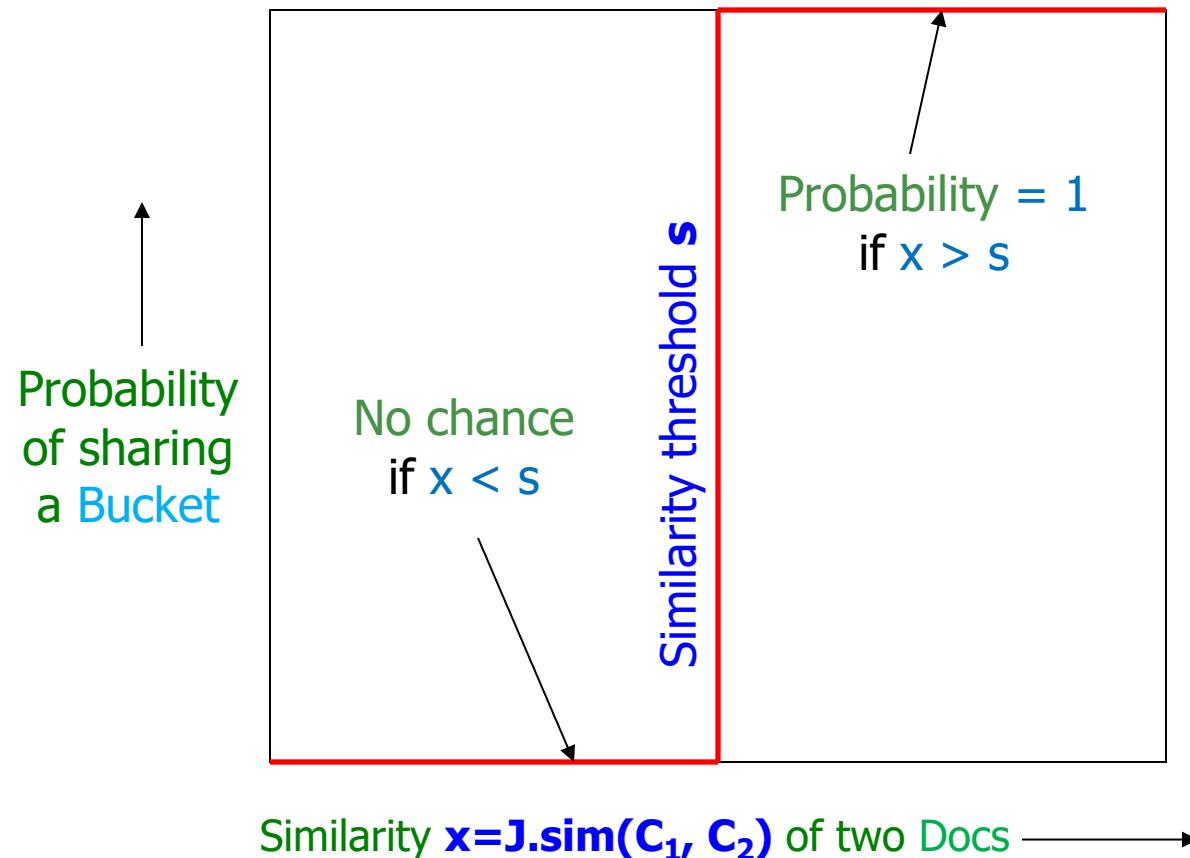
Fact 1. Probability C_1, C_2 get identical in one particular Band is
 $(0.3)^5 = 0.00243 \quad r = 5$

Fact 2. Probability C_1, C_2 get not-identical in all of the 20 Bands is $1 - (1 - 0.00243)^{20} = 0.0474$ (false positive) \rightarrow **BAD EVENT**

Confidence Rate: about 4.74% pairs of Docs with $J.\text{sim}$ 0.3%
end up becoming **candidate pairs**: They are **false positives**.
Since we will have to examine them (they are candidate pairs) but then it will turn out their similarity is below threshold s

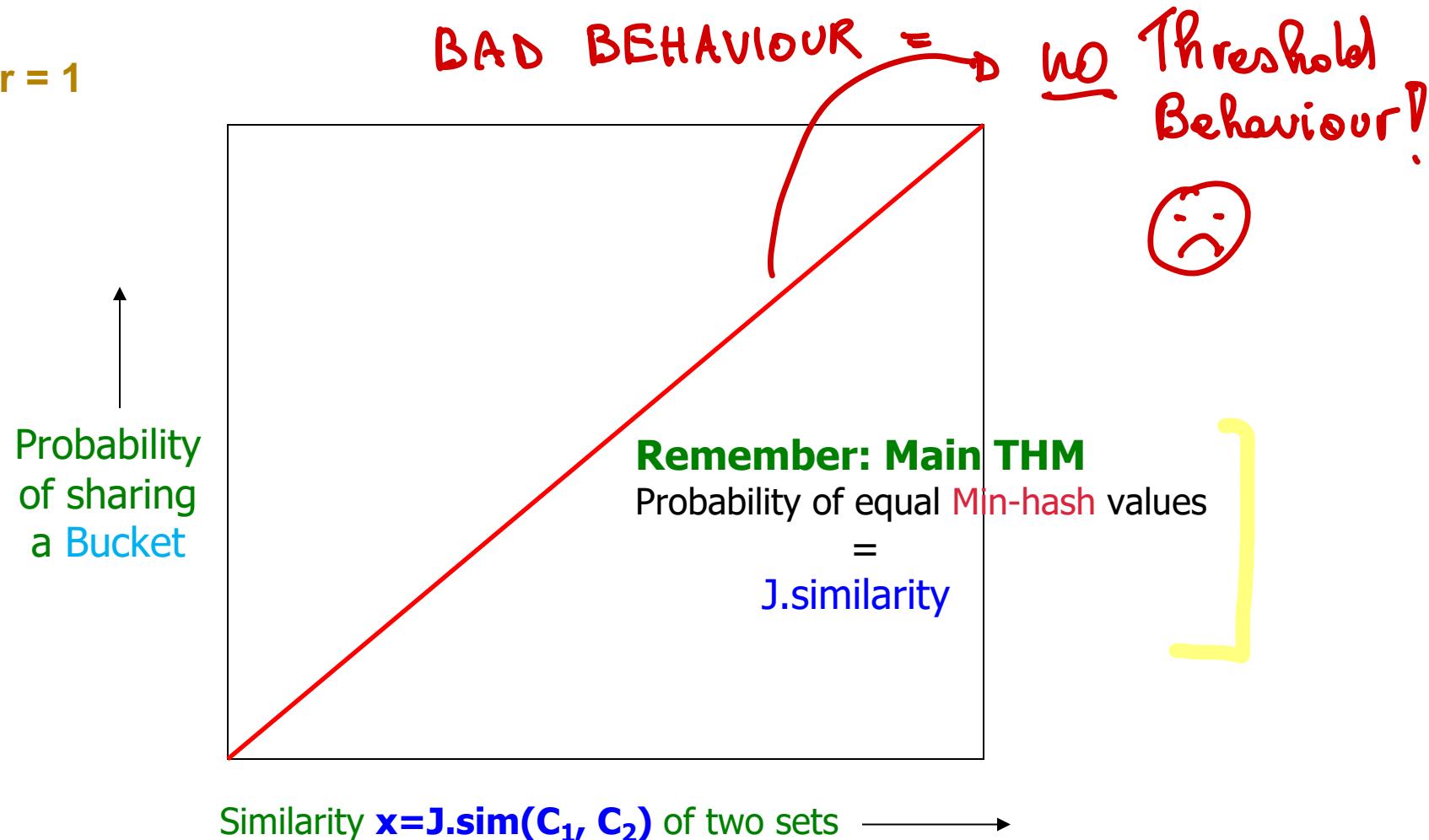
Analysis of LSH – Ideal Performance

$x = J.\text{sim}(C_1, C_2)$; s = target similarity threshold



What 1 Band of 1 Row Gives You

Case $r = 1$



Using b Bands, r rows/band

suppose Docs C_1 and C_2 have J.similarity \underline{x} .

Pick any Band (with r rows):

$$\Pr\{ \text{All rows in Band are equal} \} = x^r$$

$$\Pr\{ \exists \text{ row in Band } \text{unequal} \} = 1 - x^r$$

$$\Pr\{ \nexists \text{ Band identical} \} = (1 - x^r)^b$$

$$\Pr\{ \exists \text{ Band identical} \} = 1 - (1 - x^r)^b$$

Increases with b // (a)

Decreases with r ↘ (b)

LSH Involves a Tradeoff

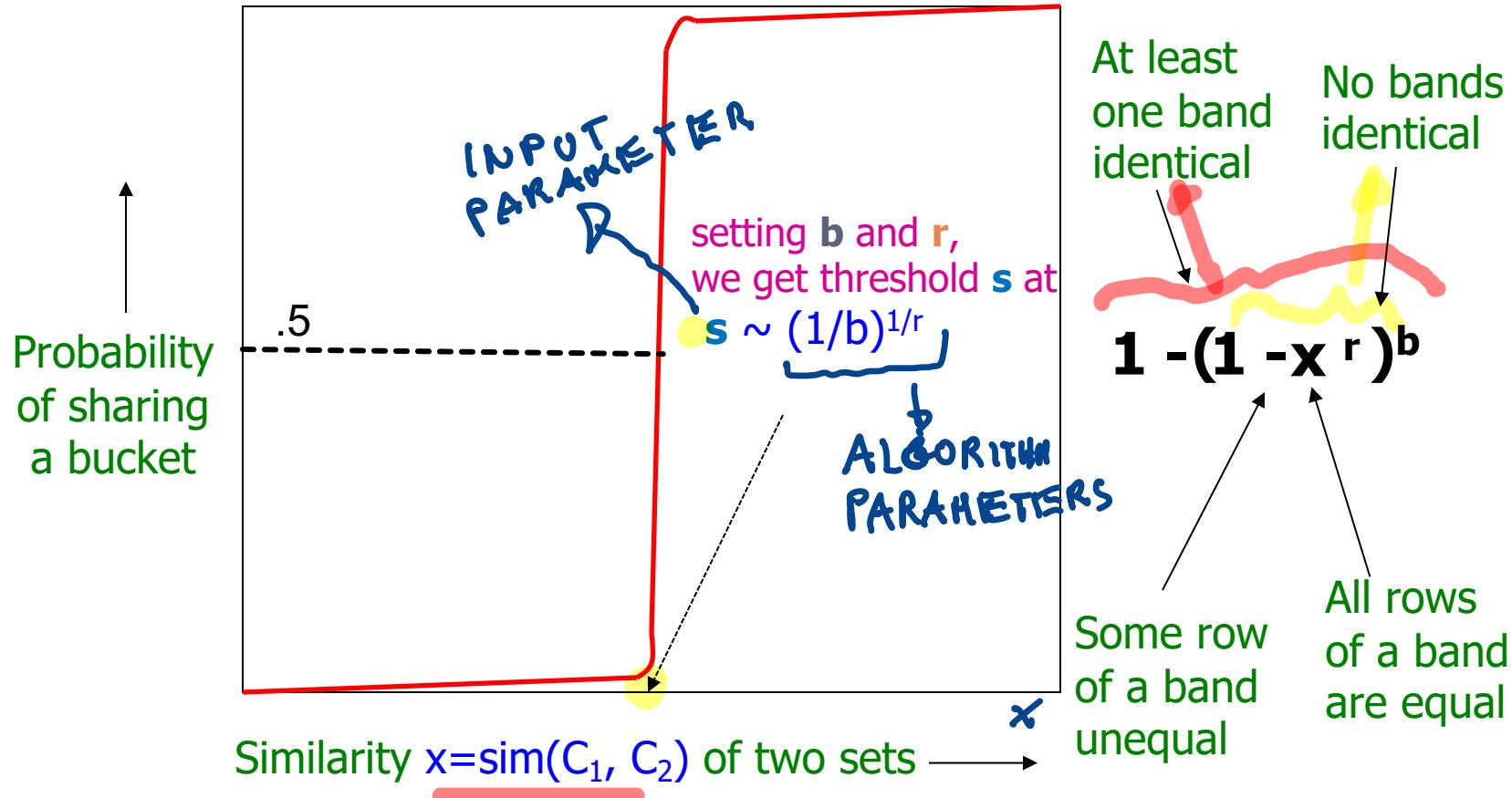
For a fixed threshold s , tune:

- The number t of Min-Hashes (rows of $SIG(*, *)$)]
 - The number of Bands b , and, hence:] (a) and (b)
 - The number of rows r per band
- to balance false positives/negatives]

Example: If we decreased to 15 Bands of 5 rows, the number of false positives (Case II before) would go down, but the number of false negatives (Case I before) would go up

What b Bands of r Rows Gives You

the case: b, r very large



LSH: The Algorithmic Goal

Task: Given the t Signature columns for each Doc (i.e. Matrix $SIG(*, *)$) and the *sim threshold* s (with $0 < s \leq 1$), **find all** the Doc pairs having: Signature Sim \geq s

LSH Solution: Apply the LSH Band technique to $SIG(*, *)$, by setting parameters b and r such that: $b \cdot r \leq t$ and $(1/b)^{1/r} \approx s$

Output: All (candidate) pairs of Docs having the same signature in all the r rows of at least one Band

Possible Parameter Tuning:

- increasing $b \rightarrow$ decreases the number of **false negatives**
- increasing $r \rightarrow$ decreases the number of **false positives**

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

$$J.Sim(C_1, C_2) = x ; b = 20; r = 5$$

Prob. that at least 1 Band is identical:

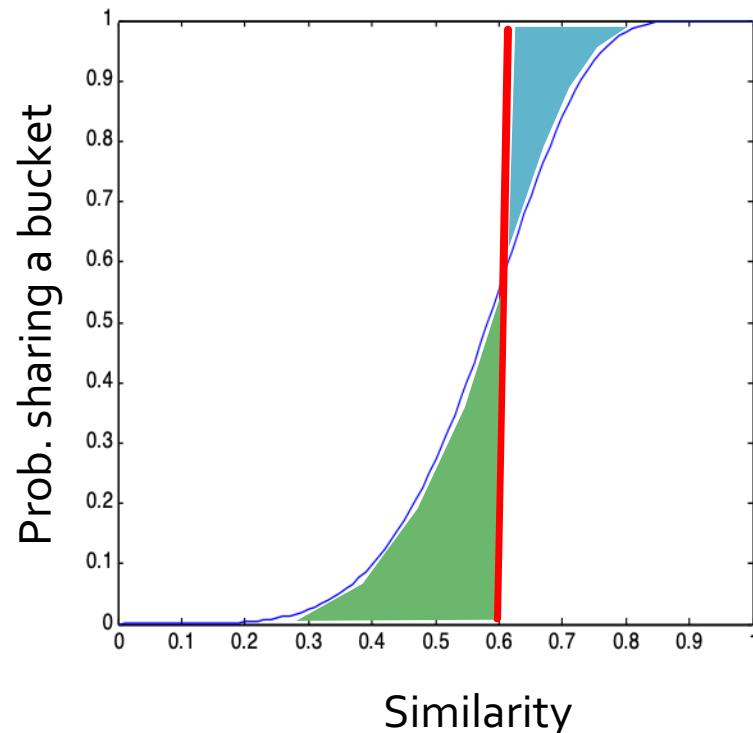
x	$1 - (1 - x^r)^b$
.2	.006
.3	.047
.4	.186
.5	.470
.6	.802
.7	.975
.8	.9996

GOOD
GAPS !

Picking r and b : The S-curve

Picking r and b to get the *best S-curve*

- $t=50$ hash-functions ($r=5$, $b=10$)



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

LSH Summary

Tune b and r to get almost all pairs with similar signatures, but eliminate most pairs that do not have similar signatures: most-balanced choice is: $s \sim (1/b)^{1/r}$

Check in main memory that **candidate pairs** really do have **similar signatures** (i.e. same values in several rows)

Optional: In another pass through data, check that the selected candidate pairs really represent similar **Docs**

↳ avoid false pos

Summary: 3 Steps

Shingling: Convert Docs to Sets

- We used **hashing** to assign each k-shingle an *ID (one row of the Big Matrix M)*

- S.SIM. for SUBSETS

Min-Hashing: Convert large sets to short signatures, while preserving *similarity*

- We used **similarity-preserving hashing** to generate signatures with property $\Pr[h_\pi(C_1) = h_\pi(C_2)] = \text{sim}(C_1, C_2)$
- We used **hashing** to get around generating random *permutations*

Locality-Sensitive Hashing: Focus on *pairs* of signatures likely to be from *similar* docs

- We used **hashing** to find candidate pairs of similarity $\geq s$

[the BAND Techniques]

the

Summary: 3 Steps

Shingling: Convert Docs to Sets

- We used **hashing** to assign each **k-shingle** an *ID (one row of the Big Matrix M)*

Min-Hashing: Convert large **sets** to short **signatures**, while preserving *similarity*

- We used **similarity-preserving hashing** to generate signatures with property $\Pr[h_\pi(C_1) = h_\pi(C_2)] = \text{sim}(C_1, C_2)$
- We used **hashing** to get around generating random **permutations**

Locality-Sensitive Hashing: Focus on *pairs* of **signatures** likely to be from *similar* **docs**

- We used **hashing** to find **candidate pairs** of similarity $\geq s$