

CSE 321 Discrete Structures

March 3rd, 2010

Lecture 22 (Supplement): LSH

Jaccard

- Jaccard similarity: $J(S, T) = |S \cap T| / |S \cup T|$
- Problem: given large collection of sets S_1, S_2, \dots, S_n , and given a threshold s , find all pairs S_i, S_j s.t. $J(S_i, S_j) > s$

Application: Collaborative Filtering

- We have n customers: $1, 2, \dots, n$
- Each customer i buys a set of items S_i
- We would like to recommend items bought by customer j , if $J(S_i, S_j) > s$

Example

S1	S2	S3	S4	S5
Toothpaste floss	Floss mouthwash	Ipod PowerBook VideoAdapter	Ipod mouthwash	Floss toothpaste mouthwash

New customer buys mouthwash, floss;
What do you recommend ?

Application: Similar Documents

- Given n documents $1, 2, \dots, n$
- Let S_i be the set of q -grams for document i
- Want to find all pairs of “similar” documents, i.e. for which $J(S_i, S_j) > s$

Example

- You work for a copyright violation detection company
- Customers: have documents 1, 2, 3,, 10^6
- Web: has pages 1, 2, 3,, 10^{11}

- Your job is to find “almost identical documents”
- What do you do ???

The Signature Method

- For each S_i , compute a signature $\text{Sig}(S_i)$ s.t.

$$J(S_i, S_j) > s \iff \text{Sig}(S_i) * \text{Sig}(S_j) \neq \text{emptyset}$$

With high probability

The Signature Method

- Step 1: compute all pairs i, j for which $\text{Sig}(S_i) * \text{Sig}(S_j) \neq \text{emptyset}$
 - This is a join operation !
- Step 2: for all such pairs, return (i, j) if $J(S_i, S_j) > s$
 - Hopefully only a few such pairs

Both false positives and false negatives are possible

The Signature Method

Will construct the signature in two steps:

1.Minhashes

2.LSH

Minhashing

- Let π be an arbitrary permutation of the domain
- For each i , let:

$$\text{mh}(S_i) = \{\text{the smallest element in } S_i \text{ according to } \pi\}$$

Example

- The entire domain is $\{a,b,c,d,e,f,g,h\}$
- The set S_i is

$$S_i = \{a,b,c,e,f\}$$

- Suppose we choose the permutation:

$$\pi = d,g,c,h,b,f,a,e$$

Then what is $mh(S_i) = ?$

Minhashing

Main property:

$$\text{Probability}(\text{mh}(S_i) = \text{mh}(S_j)) = J(S_i, S_j)$$

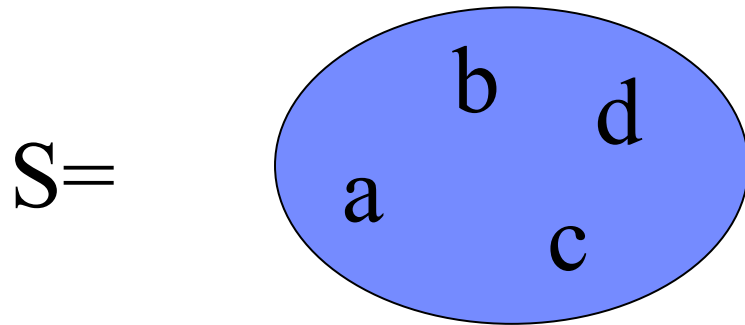
Why ?

Warmup Question

- Choose a *random* permutation π of $\{a,b,c,\dots,z\}$
- Consider the set $S=\{a,b,c,d\}$
- What is the probability that $\text{mh}(S) = c$?

Warmup Question

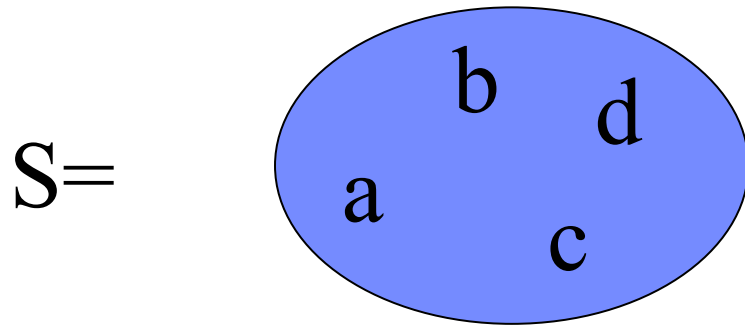
- Choose a *random* permutation π of $\{a,b,c,\dots,z\}$



What is the probability that $\text{mh}(S) = c$?

Warmup Question

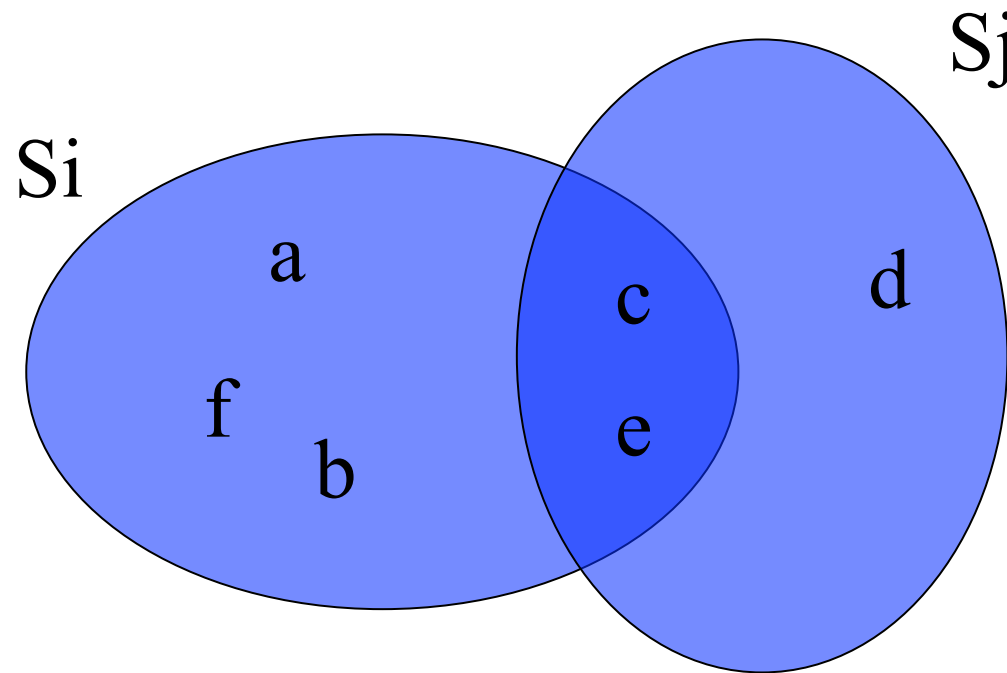
- Choose a *random* permutation π of $\{a,b,c,\dots,z\}$



What is the probability that $\text{mh}(S) = c$?

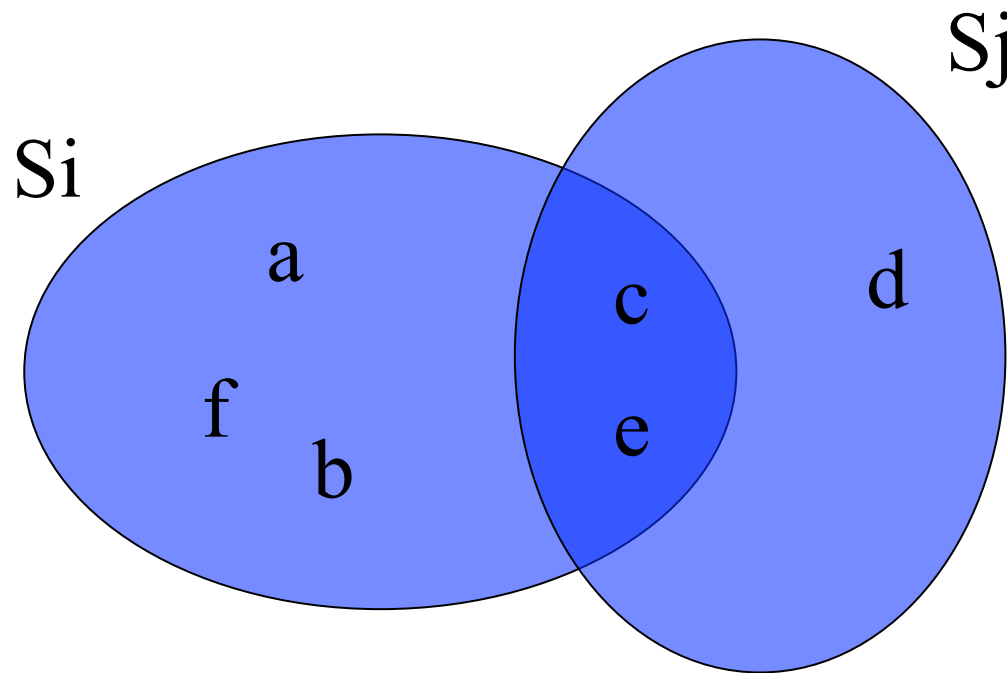
Answer: $P = \frac{1}{4}$ (each of a,b,c,d can be the min)

Main Property



What is $\text{Prob}(\text{mh}(S_i) = \text{mh}(S_j))$?

Main Property



What is $\text{Prob}(\text{mh}(S_i) = \text{mh}(S_j))$?

Computing Minhashes

- We use a hash function (which we assume is random)

```
mh(S) {  
    v = ∞;  
    forall x in S do  
        if h(x) < v then {v = h(x); y = x;}  
    return y;  
}
```

Example

- The set S_i is $S_i = \{a, b, c, e, f\}$
- Compute h :
 $h(a)=77, h(b)=55, h(c)=33, h(e)=88, h(f)=66$
- Then what is $mh(S_i) = ?$

Usage Idea

- Recall: we have n sets S_1, \dots, S_n
- Compute $mh(S_1), \dots, mh(S_n)$
- Consider only those pairs for which $mh(S_i) = mh(S_j)$: compute their Jaccard similarity
- But too many false negatives !
- How can we improve ?

Improvement

- Independent hash functions h_1, \dots, h_m
- For each S_i , compute $MH(S_i) =$ the m minhashes for each $j=1, \dots, m$

Example

- The set S_i is $S_i = \{a,b,c,e,f\}$
- Compute h_1 :
 $h_1(a)=77, h_1(b)=55, h_1(c)=33, h_1(e)=88, h_1(f)=66$
- Compute h_2 :
 $h_2(a)=22, h_2(b)=66, h_2(c)=55, h_2(e)=11, h_2(f)=44$
- Then what is $MH(S_i) = ?$

Example

- The set S_i is $S_i = \{a, b, c, e, f\}$
- Compute h_1 :
 $h_1(a)=77, h_1(b)=55, h_1(c)=33, h_1(e)=88, h_1(f)=66$
- Compute h_2 :
 $h_2(a)=22, h_2(b)=66, h_2(c)=55, h_2(e)=11, h_2(f)=44$
- Then what is $MH(S_i) = ?$
- Answer: $MH(S_i) = (c, e)$ (an ordered pair)

Using Minhashes

- Compute $MH(S_1), \dots, MH(S_n)$
- $J(S_i, S_j) \approx$ the fraction of positions where $MH(S_i)$ and $MH(S_j)$ agree

Example

n=7

S1	S2	S3	S4	S5	S6	S7
----	----	----	----	----	----	----

m=6

1	1	4	3	5	4	4
4	6	4	6	9	4	4
8	7	8	9	5	6	8
4	8	7	8	4	3	7
8	7	8	9	5	6	8
6	8	6	8	4	5	6

Estimate $J(S1,S2) = ?$, then $J(S1,S3) = ?$

Note:

Minhashes still require Jaccard

- How do we compute the fraction of positions where MH_i and MH_j agree ?
- Annotate each position with its position number
- Then this is precisely $J(MH_i, MH_j)$

Note:

Minhashes still require Jaccard

1	1		4
2	4 ←		4
3	8 ←		8
..	4		7
..	8 ←		8
m	6 ←		6

$$\text{MH1} = \{(1,1), (2,4), (3,8), (4,4), (5,8), (6,6)\}$$

$$\text{MH2} = \{(1,4), (2,4), (3,8), (4,7), (5,8), (6,6)\}$$

$$J(\text{MH1}, \text{MH2}) = 4 / (2m - 4)$$

Fraction of equal positions = $4/m$

Comments on Minhashes

- It is not a signature yet !
- We have only reduced the problem of computing $J(S_i, S_j)$ to the problem of computing J on smaller sets, of size m
- The signature is provided by LSH (next)

Locality Sensitive Hashing

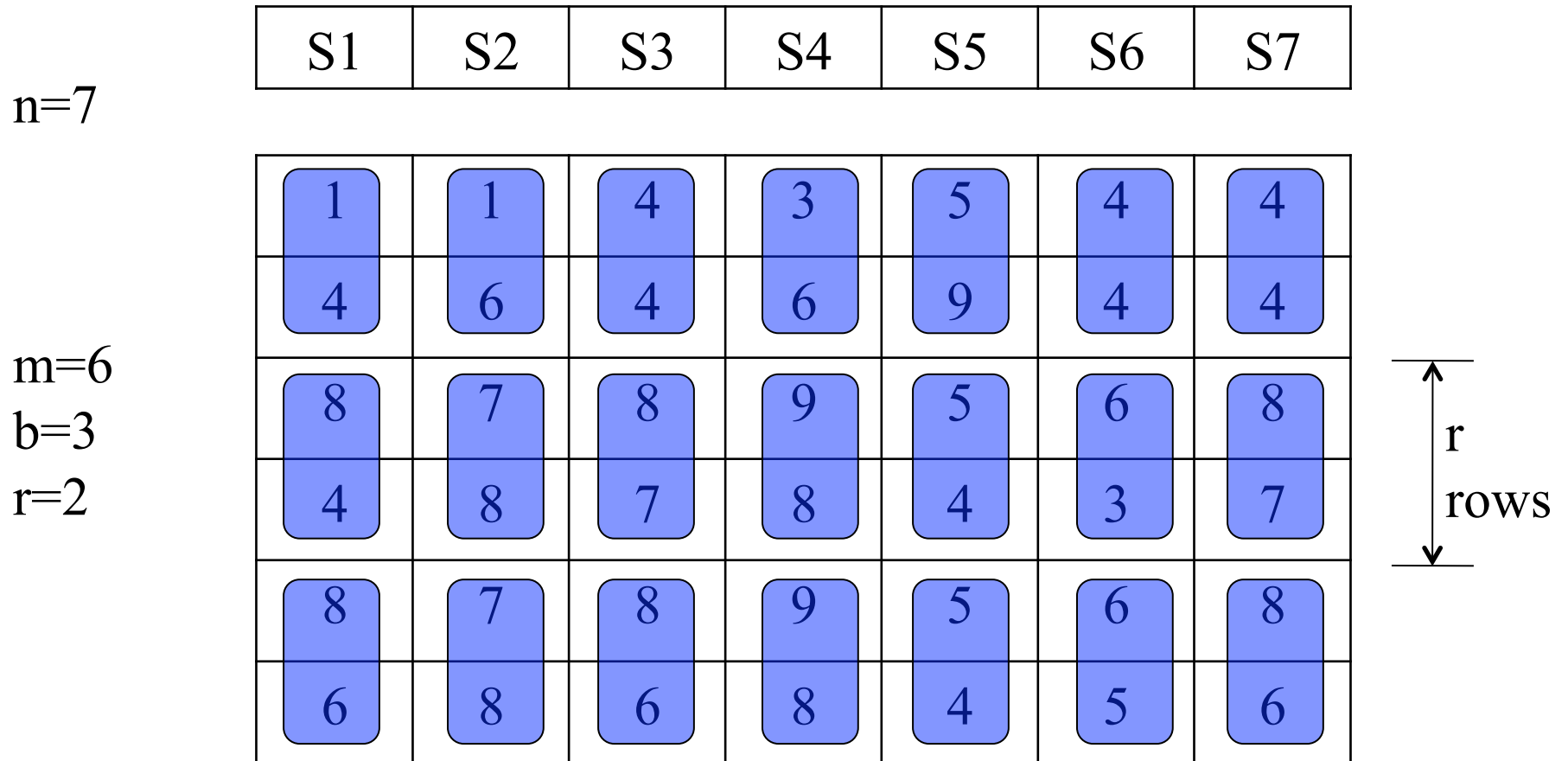
- We have n strings MH_1, \dots, MH_n , each of length m
- Compute a signature $\text{Sig}(MH_i)$ for each string MH_i
- Desired property:

$$\text{Sig}(MH_i) * \text{Sig}(MH_j) \neq \text{empty} \iff J(MH_i, MH_j) > s$$

LSH

- Each MHi consists of m values
- Each signature $Sig(MHi)$ will have size b
- Divide the m values into b “bands” of size r “rows”, i.e. $m = b * r$
- For each band $j=1, \dots, r$, apply a hash function h_j to the string of values in band j in $MHi \rightarrow h_j$
- Then $Sig(MHi) = (h_1, h_2, \dots, h_b)$

Example



Analysis

- Goal: want to compute the probability that $\text{Sig}(\text{MH}(S_i)) * \text{Sig}(\text{MH}(S_j)) \neq \text{emptyset}$, as a function of $s = J(S_i, S_j)$
- So let $s = J(S_i, S_j) = \text{fraction of equal positions}$

Analysis

$$J(S_i, S_j) = s$$

What is the probability that two entries are equal ?

S_i	S_j
7	7

Analysis

$$J(S_i, S_j) = s$$

What is the probability that two entries are equal ?

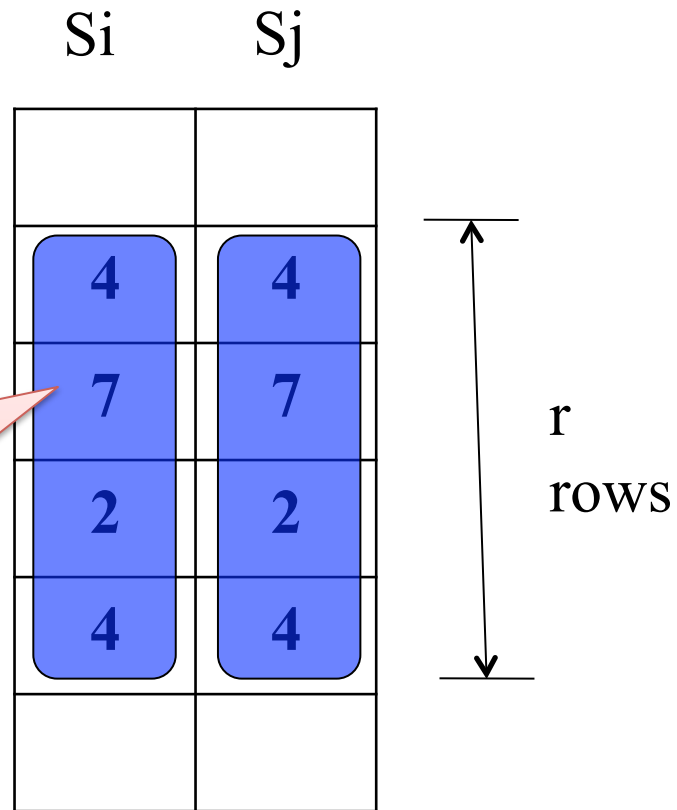
S_i	S_j
7	7

Answer: s

Analysis

$$J(S_i, S_j) = s$$

What is the probability that the bands are equal ?



Analysis

$$J(S_i, S_j) = s$$

What is the probability that the bands are equal ?

S_i	S_j
4	4
7	7
2	2
4	4

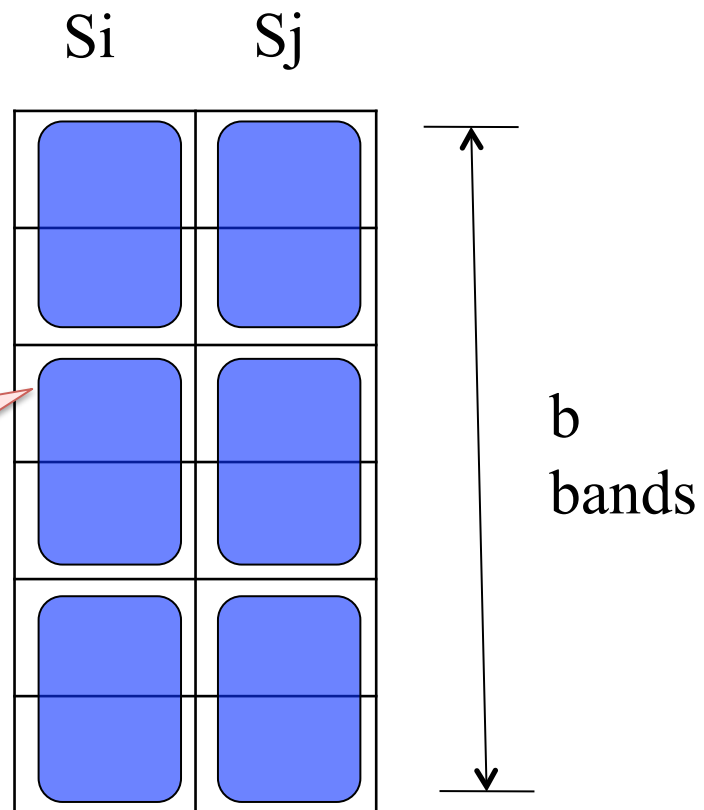
r
rows

Answer: s^r

Analysis

$$J(S_i, S_j) = s$$

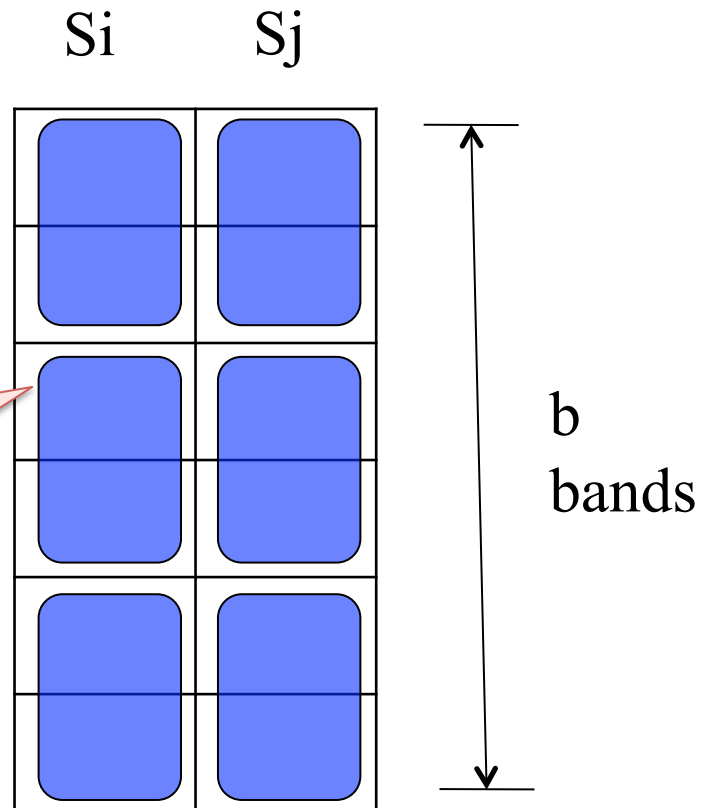
What is the probability that some pair of bands are equal ?



Analysis

$$J(S_i, S_j) = s$$

What is the probability that some pair of bands are equal ?

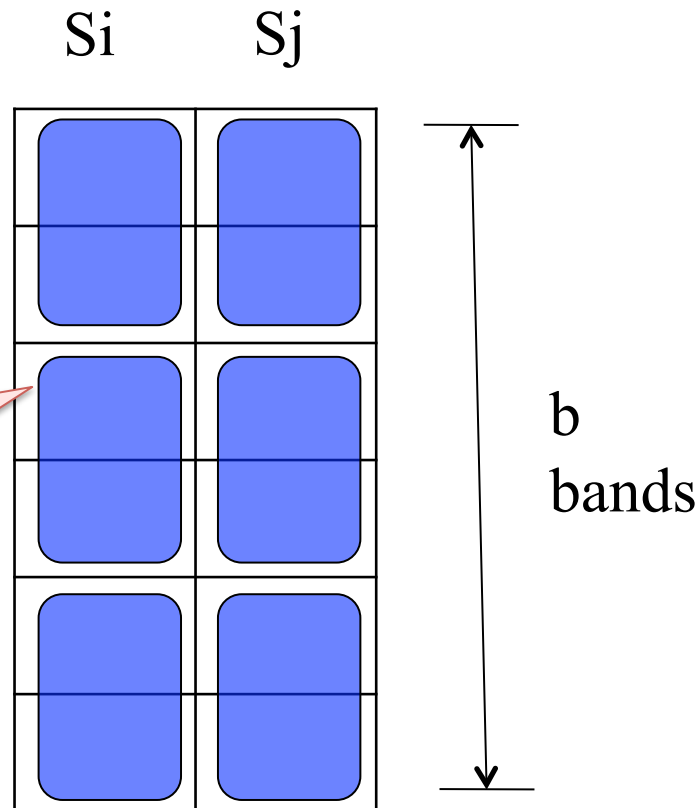


Answer: $1 - (1-s)^b$

This is precisely
the probability that
 $\text{Sig}(S_i) * \text{Sig}(S_j) \neq \text{emptyset}$

$$J(S_i, S_j) = s$$

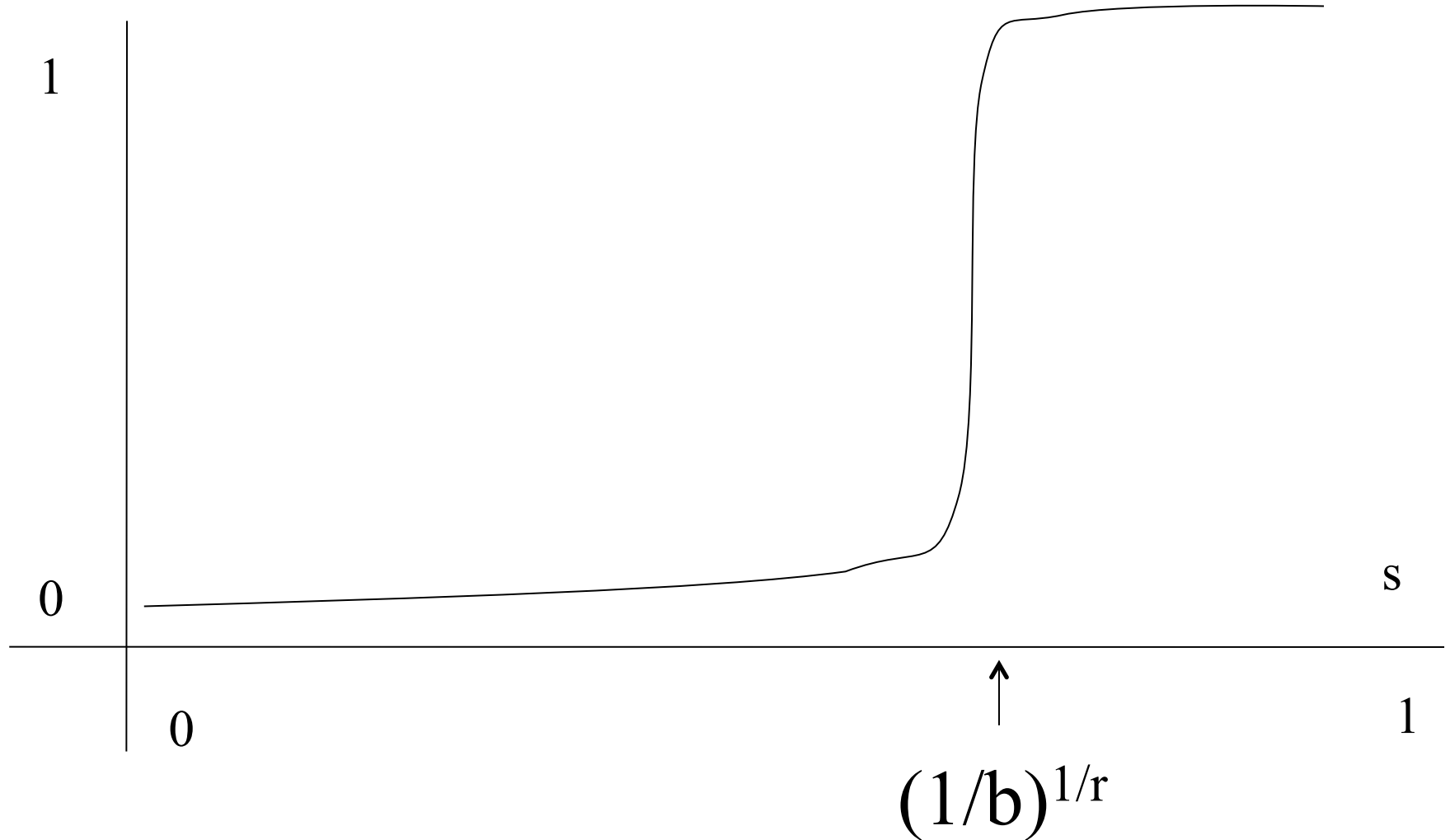
What is the
probability
that some pair of
bands are equal ?



Answer: $1 - (1 - s^r)^b$

Analysis

Probability $1-(1-s^r)^b$



Putting it together

- You work for a copyright violation detection company
- Customers: have documents 1, 2, 3,, 10^6
- Web: has pages 1, 2, 3,, 10^{11}

- Your job is to find “almost identical documents”
- What do you do ???

Step 1: Q-grams

Doc

DocID	Qgram
1	To
1	o b
1	be
1	be
1	e o
1	or
1	or
1	r n
1	no
...	...
...	
2	..
...	

Web

URL	Qgram
abc.com	Oba
abc.com	Bam
abc.com	ama
abc.com	ma
abc.com	a h
...	...
bcd.com	...
...	

Step 2: Compute m Min-hashes

docID	mh1	mh2	mh3	mh4	...	mh500
1	To	que	Ion	or		Not
2				

url	mh1	mh2	mh3	mh4	...	mh500
abc.com	sec	ret	def	ens		...
bcd.com				

Note: $m =$ a few hundreds

Step 3: Compute Signatures

docID	$h(\text{mh}_1, \dots, \text{mh}_{20})$	$h(\text{mh}_{21}, \dots, \text{mh}_{40})$...
1	2345234	3232	...
2	...		

docSIG

docID	Sig
1	2345234@1
1	3232@2
1	...
1	452342@25
2	23423@1
2	...

webSIG

url	Sig
abc.com	676876@1
abc.com	3232@2
abc.com	...
abc.com	787892@25
bcd.com	23423@1
...	...

Step 4: Find docs with common signatures

Need two indexes:

for every signature s , $\text{doc}[s]$ = set of documents containing s

for every webpage w , $\text{web}[]$ = set of webpages containing s

```
For all  $s$  in Sig do
  for  $d$  in  $\text{doc}[s]$ , for  $w$  in  $\text{web}[s]$  do
    print( $d, w$ )
```