

Week 6

min Hash

represents each document as a small number of numbers from which we can easily estimate the Jaccard Distance

N documents/string
(d average size)

↓ Preprocess them

Given a query document
report if any document
stored is close to it

$m = \#$
of queries

LSH "Locality - Sensitive hashing"
This allows efficient search for things that are similar

"Jaccard distance"

Break each document into fragments (size 7-8)
typical

$$d(x, y) = \frac{\text{fragments in } x \text{ and } y}{\text{fragments in } x \text{ or } y}$$

n documents (average size d)

m queries

$$O(m \cdot n \cdot d)$$

get rid of this
I

get rid of this
II

min Hash

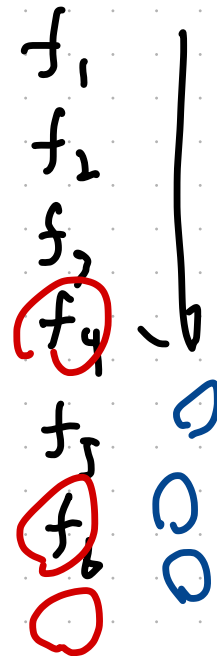
Have n documents,

F = the set of all their fragments
in random order

$\text{Score}(x)$ = The index of the first
document a fragment of x
fragment in F that is

$\Pr[\text{score}(x) = \text{score}(y)] = \text{Jaccard distance}$
between x and y

$$E \left[\begin{array}{l} 1 \text{ if } \text{score}(x) = \text{score}(y) \\ 0 \text{ if } \text{score}(x) \neq \text{score}(y) \end{array} \right] = \text{J}$$



LSH: Locality Sensitive Hashing

- Suppose you have a set H of hash functions
- Same data
- Same distance function on the data

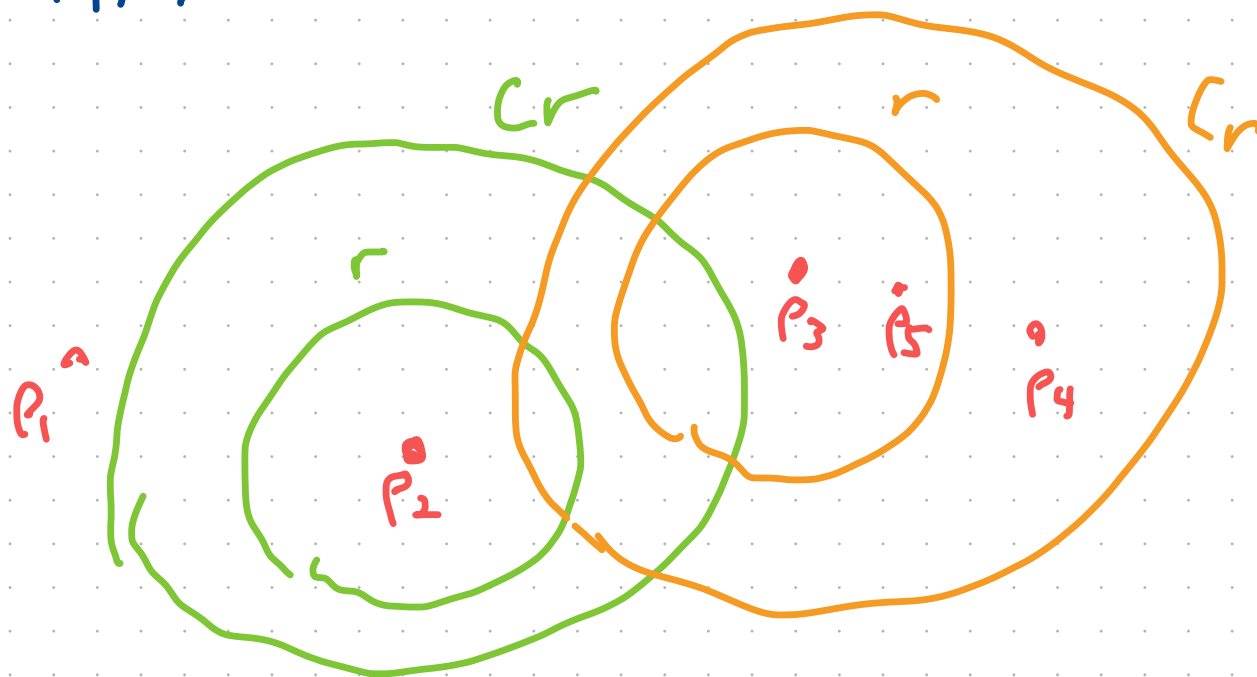
- If $d(x, y) \leq r$ then $h(x) = h(y)$ with prob $\geq p_1$
- If $d(x, y) \geq cr$ then $h(x) = h(y)$ with prob $\leq p_2$

$$p_1 \geq p_2$$

want $p_1 = 1$

$p_2 = 0$

$$\text{Error} = \min(1 - p_1, p_2)$$



Hamming Distance

s_1	a	a	a	b	b	b	c	c	c
s_2	a	a	c	b	b	a	c	c	d
	0	1	2	3	4	5	6	7	8

$$d_{\text{Hamming}}(s_1, s_2) = 3$$

Pick a random $\text{int}(i)$ in the range $1..N$

Let $h(s) = s[i]$

$$i = 6$$

$$h(s_1) = c$$

$$h(s_2) = c$$

LSH: Locality Sensitive Hashing

- Suppose you have a set H of hash functions
- Same data
- Same distance function on the data

• If $d(x, y) \leq r$ then $h(x) = h(y)$ with prob $\leq p_1$

• If $d(x, y) \geq cr$ then $h(x) = h(y)$ with prob $\geq p_2$

$p_1 > p_2$

hamming

Choose

$$r = 0.8n$$

$$cr = 0.9n$$

$$0.8 = p_1$$

$$0.9 = p_2$$

- Q.8 P_1 = Chance that distance $\leq r$, correct 0.8
- Q.9 P_2 = Chance that distance $\geq cr$, incorrect 0.9
- answer

"AND" construction.

Hash k times, say distance small
if all k are the same

$$P_1 \rightarrow P_1^k$$

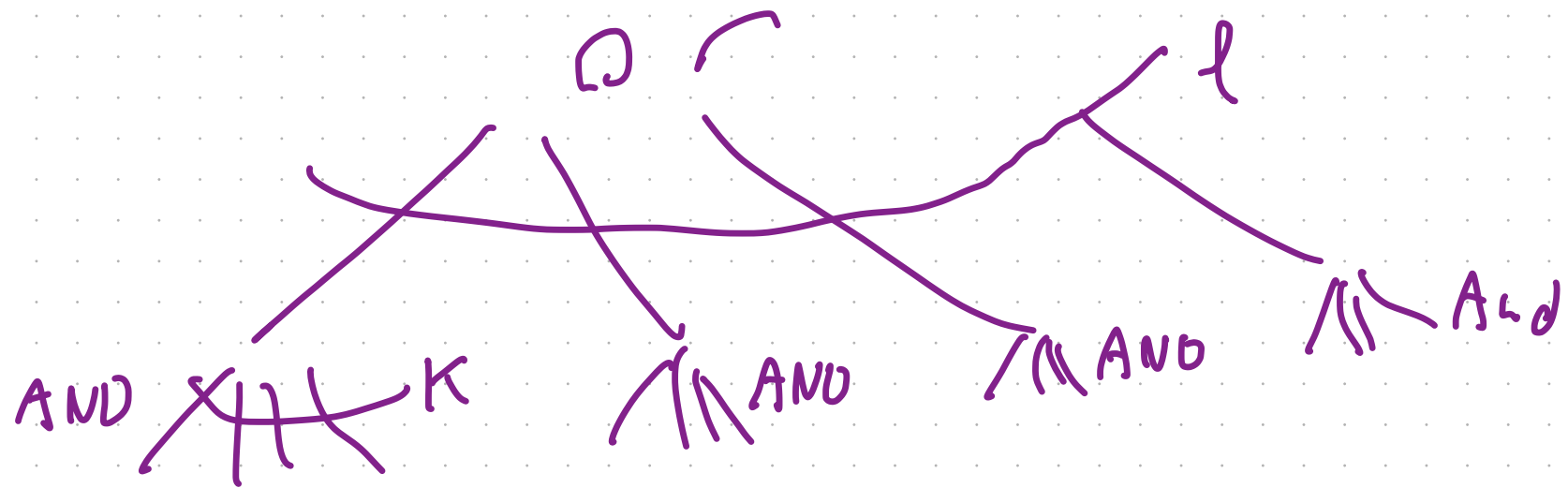
$$P_2 \rightarrow P_2^k$$

"OR" construction

Hash l times
Distance is small if
at least 1 is small

$$P_1 \rightarrow 1 - (1 - P_1)^k$$

$$P_2 \rightarrow 1 - (1 - P_2)^k$$



$$P_1 \rightarrow 1 - (1 - P_1^K)^l$$

$$P_2 \rightarrow 1 - (1 - P_2^K)^l$$

Data



Fragments



Multiple Scores ($l \times K$)



Group them into l groups of K



for each group, make a dictionary that maps the scores of the group to the string that made it.

To search for s , compute scores, group, for each group look in the dictionary

$$D[0] = \begin{cases} [1, 5] \rightarrow s_1 \\ [1, 7] \rightarrow s_2 \end{cases}$$

$$D[1] = \begin{cases} [2, 3] \rightarrow s_1, s_2 \end{cases}$$

