

Hashing  $\rightarrow$  64-bit hash  $\rightarrow$  ?

Duplicate / near-duplicates

$\approx 30\%$  of web pages

Dealing with near-duplicates

Sets of #'s

2 sets A, B, let each be of ~~size~~  $n$

$\in$  64-bit #'s

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

1 if  $A = B$

0 if disjoint

if ordered  $O(n)$

$O(n \log n)$  for ordering

unordered  $O(n^2)$

} too expensive for each document comparison

$\downarrow$

change problem

$\downarrow$

approximate resemblance

assume

A black box for random permutation

$$\pi_3(x) \rightarrow y$$

x	0	1	2	3	4	5	6	7
$\pi_3(x)$	2	4	3	1	7	6	5	0

$$\pi_3(2) = 3$$

permutation of a set, function applied to each element

$$A = \{2, 4, 6\} \quad \pi_3(A) = \{3, 5, 7\} \quad \min \pi_3(A) = 3$$

$$\min(\pi_1(A)), \min(\pi_2(A)), \min(\pi_3(A)) \dots \min(\pi_{100}(A))$$

keep 100 ~~pts~~ associated with each set ~~across~~ across all sets  
 $\rightarrow$  calling card for each set

A calling card

$$\min(\pi_1(A)), \min(\pi_2(A)), \min(\pi_3(A)) \dots$$

B calling card

$$\min(\pi_1(B)), \min(\pi_2(B)), \min(\pi_3(B)) \dots$$

$$\text{Prob}[\min \pi_1(A) = \min \pi_1(B)] = \text{Resemblance}(A, B)$$



$$A = \{2, 4, 6\} \quad \pi_3(A) = \{3, 5, 7\} \quad \min \pi_3(A) = 3$$

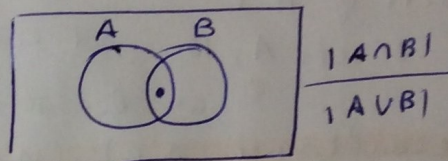
$$B = \{4, 5, 6\} \quad \pi_3(B) = \{5, 6, 7\} \quad \min \pi_3(B) = 5$$

every number has an <sup>equal</sup> chance to be ~~the~~ the minimum of a permutation function. If ~~a~~ a number is present in both documents (sets) and one of the "right" permutations is applied, the minima match

→ an elt can be min in several  $\pi$ s

→ when minima match an elt is in both sets

→ given a doc every elt has equal chance to be min, ~~these~~ random



go through 100  $\pi$ s  
→ get an estimate of resemblance

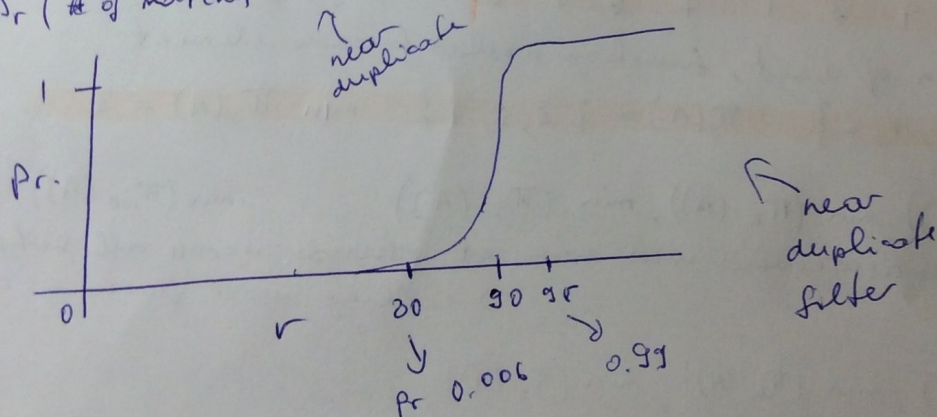
Suppose true resemblance is  $r$

$\geq 90$  matches → similar

$< 90$  not

$$\Pr(\# \text{ of matches} = 90 \text{ when resemblance} = r) = \binom{100}{90} r^{90} (1-r)^{10}$$

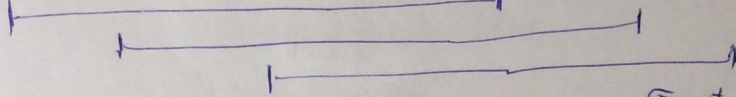
$$\Pr(\# \text{ of matches} \geq 90 \text{ when resemblance} = r) = \sum_{k=90}^{100} \binom{100}{k} r^k (1-r)^{100-k}$$



(2)

Turn a document into a set of #s  
x ↗ semantic

Four score and seven years ago. —



$k=4$   
then hash

shingling:

Syntactic vs. semantic

↖ takes some context  
into account