

## SIMILARITY HASH

MIN-HASH = minwise independent permutations

Permutations  $X = \text{set of items}$   $|X|!$  permutations

$$X = \{a, b, c\}$$

$$\begin{matrix} a & b & c \\ \cancel{a} & \cancel{b} & \cancel{c} \\ a & c & b \end{matrix}$$

$$\begin{matrix} b & a & c \\ \cancel{b} & \cancel{a} & \cancel{c} \\ b & c & a \end{matrix}$$

$$\begin{matrix} c & a & b \\ \cancel{c} & \cancel{a} & \cancel{b} \\ c & b & a \end{matrix}$$

$3!$  permutations to

Each permutation induces an order (e.g.  $cab \rightarrow c < a < b$ )

We take the minimum of all permutations

$h(X) = \text{permutation of } X$

$$\min h(X) \stackrel{\text{def}}{=} \underset{a \in X}{\text{arg-min}} h(a)$$

$\min h(X) : X = \cup \rightarrow a \in X$

example:

$$\begin{array}{c|ccc} & c & a & b \\ \hline h() & | & 1 & 2 & 3 \end{array}$$

$$\min h(X) = c$$

Given a set  $X \subseteq U$ , permutation  $h: U \rightarrow [|U|]$  (ranking)

family  $\mathcal{H} = \{h: X \rightarrow [|X|] \text{ bijective}\}$  is MIN-WISE INDEPENDENT  
if  $\forall X \subseteq U \quad \forall a \in X$

$$\Pr_{h \in \mathcal{H}} [a = \min h(X)] = \frac{1}{|X|}$$

We like this because of Jaccard's index for similarity  
(see Broder & Altunis)

document = bag of words = set of "lemmatized" word

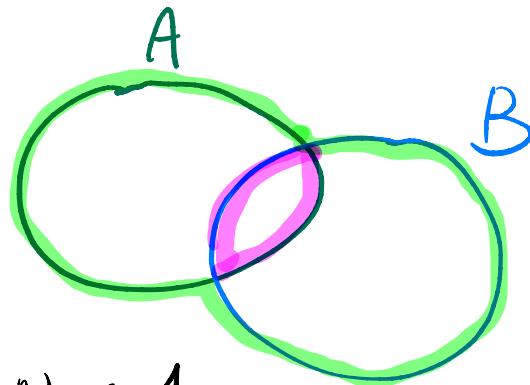
$U = \{\text{all possible words}\}$

$A \subseteq U$  (documents)

$A, B \subseteq U$

Jaccard's index

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



$$0 \leq J(A, B) \leq 1$$

$$A \cap B = \emptyset$$

$$A = B$$

After sorting  $A, B$  :  $J(A, B)$  takes  $\mathcal{O}(|A| + |B|)$  time

Big data :  $J(A, B)$  is too expensive for zillions of  $A, B$

MIN-HASH  
(SIM-HASH)

Interesting property:  $\Pr_{h \in \mathcal{H}} (\min h(A) = \min h(B)) = J(A, B)$   $\forall A, B$

TOO SEE WHY

$$r = |A \setminus B|$$

$$b = |A \cap B|$$

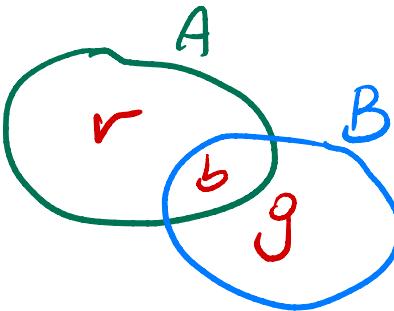
$$g = |B \setminus A|$$

$$|A| = r + b$$

$$|B| = b + g$$

$$|A \cap B| = b$$

$$|A \cup B| = r + b + g$$



$$\Pr(\min h(A) = \min h(B)) = \Pr[\exists y \in A \cap B : y = \min h(x)]$$

$$= \sum_{y \in A \cap B} \underbrace{\Pr[y = \min h(x)]}_{\frac{1}{|X|}} = \frac{b}{r+b+g} = \frac{|A \cap B|}{|A \cup B|} = J(A, B)$$

$$X = A \cup B$$

$$|X| = r + b + g$$

MIN-HASH

$$\Pr[\alpha = \min h(x)] = \frac{1}{|X|}$$

## Some observations

- Storing  $h(U)$  is unfeasible in many applications
- One single  $h \in \mathcal{H}$  could not give enough precision

Fix

- It has been shown that we can "safely" replace permutation family  $\mathcal{H}$  with our UNIVERSAL family  $\mathcal{H}$
- Use  $k$  hash functions  $h_1, h_2, \dots, h_k$

$k$ -min hash  $A \rightarrow \langle \min h_1(A), \min h_2(A), \dots, \min h_k(A) \rangle = S(A)$  sketch

$$\text{approximate Jaccard : } \tilde{j}_{(A,B)} = \frac{|S(A) \cap S(B) \cap S(A \cup B)|}{|S(A \cup B)|}$$

bottom- $k$  hash:

instead of using  $k$  hash functions  
we use a single  $h$ , and take  
the  $k$  smallest ranked items

$$S(A \cup B) = \{a, b, c, d, f, g\}^{k=4}$$
$$S(S(A) \cup S(B)) = \{a, b, f, g\}$$
$$\{a, b, c, d\}$$
$$S(A) = \{a, c, d, f\}$$
$$S(B) = \{a, b, f, g\}$$
$$a < b < c < d < e < f < g$$

Given  $S(A), S(B)$ , then  $J(A, B)$  can be computed in  $O(k)$  time

$\Rightarrow$  Step 1 compute  $S(A)$  for each set  $A$

Step 2 for each pair  $A, B$  : compute  $\tilde{J}(A, B)$

much faster than computing  $J(A, B)$

Price to pay: it's an approximation. Use indicator variables

$k$ -min hash : given  $A, B$

$$X_i = \begin{cases} 1 & \text{if } \min h_i(A) = \min h_i(B) \\ 0 & \text{otherwise} \end{cases} \quad pr = J(A, B)$$

$$E[X_i] = J(A, B) \quad \text{UNBIASED ESTIMATOR}$$

$$Y = \frac{\sum_{i=1}^k X_i}{K} \quad \text{UNBIASED ESTIMATOR} \quad E[Y] = \frac{\sum_{i=1}^k E[X_i]}{K} = \frac{k J(A, B)}{K} = J(A, B)$$

$$Y = \tilde{J}(A, B)$$

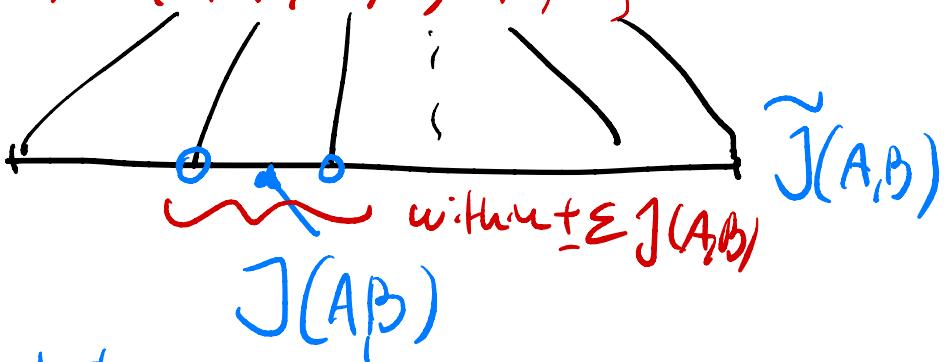
$$k=1 \rightarrow Y \in \{0, 1\}$$

$$k=2 \rightarrow Y \in \{0, \frac{1}{2}, 1\}$$

$$k=3 \rightarrow Y \in \{0, \frac{1}{3}, \frac{2}{3}, 1\}$$

$$k=4 \rightarrow Y \in \{0, \frac{1}{4}, \frac{2}{4}, \frac{3}{4}, 1\}$$

$$k \dots \rightarrow Y \in \{0, \frac{1}{k}, \frac{2}{k}, \dots, \frac{k-1}{k}, 1\}$$



Chernoff's bound helps us to evaluate

CHERNOFF'

$$Y' = \sum_{i=1}^k X_i \quad (Y' = kY) \Rightarrow \mu' = E[Y'] = k \underbrace{E[Y]}_{\mu} = k J(A, B)$$

$$\Pr[|Y - \mu| \geq \epsilon \mu] \stackrel{\text{multiply by } k}{=} \Pr[|Y' - \mu'| \geq \epsilon \mu'] \stackrel{\text{CB}}{\leq} 2 \cdot e^{-\frac{\epsilon^2 \mu'}{3}} = 2e^{-\frac{\epsilon^2 k \mu}{3}}$$

relative error

Is this good? We want to bound this probability by  $\delta < 1$

$$2e^{-\frac{\epsilon^2 k \mu}{3}} < \delta \Rightarrow k = O(\epsilon^{-2} \lg \delta^{-1} \mu^{-1})$$

not good  $\nabla \mu = J(A, B)$

$$\mu^{-1} \sim |A| H(B)$$

Let's use a more suitable version of CBs (concentration bounds)

AZUMA-Hoeffding bound (AH)

$X_1, X_2, \dots, X_K$  i.i.d. random vars,  $\mu = E\left[\frac{\sum_{i=1}^K X_i}{K}\right]$

$$a_i \leq X_i \leq b_i$$

$\hookrightarrow$  our  $X_i \in \{0, 1\}$

$$\Pr[|Y - \mu| \geq \varepsilon] \leq 2 \cdot e^{-\frac{2K^2\varepsilon^2}{\sum_{i=1}^K (b_i - a_i)^2}}$$

no  $\mu$  here      ABSOLUTE ERROR  $\Theta$

IN OUR CASE:

$$2 \cdot e^{-\frac{2K^2\varepsilon^2}{K}} = 2e^{-2n\varepsilon^2} < \delta \Rightarrow K = O(\varepsilon^{-2} \lg \delta^{-1})$$

( $\because$ )

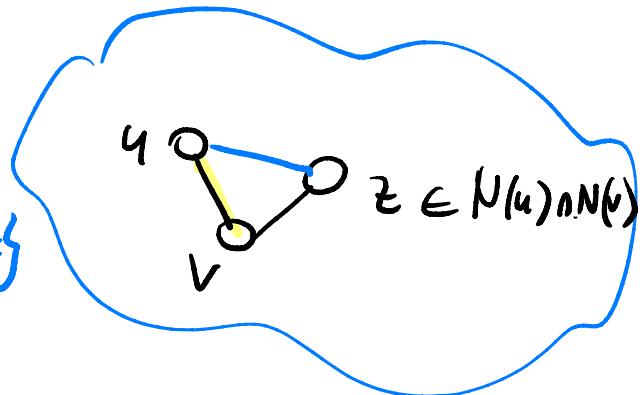
$$a_i = 0 \quad b_i = 1 \Rightarrow \sum_{i=1}^K (b_i - a_i)^2 = K$$

NETWORK ANALYSIS : social networks "triangle" is a measure of sociability

$$G = (V, E)$$

$$x \in V$$

$$N(x) = \{y : xy \in E\}$$



$G$

Counting # triangles

BASIC ALGORITHM :  
for each edge  $uv \in E$ :

$$C += |N(u) \cap N(v)|$$

$$\text{return } C/3$$

$O(|E|^3)$

Let's estimate instead  $C = \# \text{triangles}$

we use Jaccard and min-hash

$$A = N(u) \quad B = N(v)$$

$$|A \cap B|$$

$$\text{J}(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

....

$$|A \cap B| = \text{J}(A, B) \cdot (|A| + |B|)$$

Using min-hash, we get an estimation

$$|\tilde{A \cap B}| = \frac{\tilde{\text{J}}(A, B) (|A| + |B|)}{1 + \tilde{\text{J}}(A, B)}$$

note: error is both at the numerator  
and denominator

TIME IS LINEAR  $O(|E|)$

