

# Min Hashing

Note Title

1/25/2016

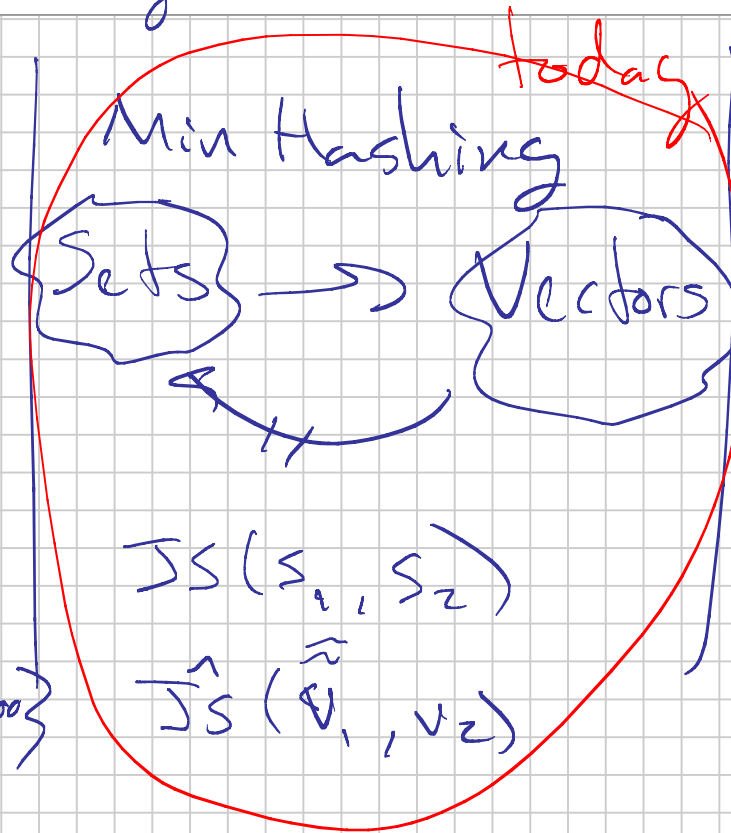
K-grams

Documents



Sets

$\{1, 2, 7\} \subset \{1, 2, \dots, 100\}$



LSH

• generalize

$X = \text{large \# sets}$

+  $g \leftarrow \text{binary}$

if  $g \in X$

+ are any two  $g_i, g_j \in X$  same

Set  $\rightarrow$  vector  $v = (v_1, v_2, \dots, v_k)$  *i.i.d*

$$JS(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|23|}{|\{1, 2, 3, 4\}|} = \frac{2}{4} = \frac{1}{2}$$

$A = \{1, 2\}$   
 $B = \{2, 3, 4\}$

# Matrix

	$S_1$	$S_2$	$S_3$	$S_4$
1	1	0	0	1
2	1	0	1	0
3	0	1	1	0
4	0	0	1	1
5	1	0	0	0
6	0	0	1	1

$$S_1 = \{1, 2, 5\}$$

$$S_2 = \{3\}$$

$$S_3 = \{2, 3, 4, 6\}$$

$$S_4 = \{1, 4, 6\}$$

← very sparse

## Step 1 Permute Rows Step 2

	$S_1$	$S_2$	$S_3$	$S_4$
2	1	0	1	0
3	1	0	0	0
6	0	0	0	1
4	0	0	1	0
5	1	0	0	0

$$m(S_1) = 2$$

$$m(S_2) = 3$$

$$m(S_3) = 2$$

$$m(S_4) = 6$$

$$\hat{J}S_i = (S_i, S_j) = \begin{cases} 1 & \text{if } m(S_i) = m(S_j) \\ 0 & \text{otherwise} \end{cases}$$

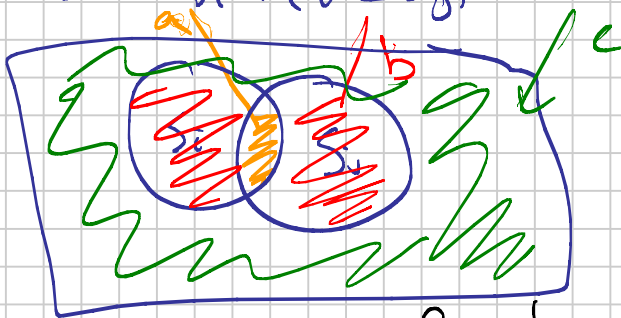
$$E[\hat{J}S_i(S_i, S_j)] = JS(S_i, S_j)$$

$$= |S_i \cap S_j|$$

$$|S_i \cup S_j| = |S_i \cap S_j| + |S_i \Delta S_j|$$

$$X_0 = S_0 \cup S_6$$

$$= |S_0 \cap S_6| + |S_0 \Delta S_6|$$



$$\Pr(m_i = m_j) = \frac{a}{a+b}$$

first non-zero

$$\hat{JSS}_k(A, B) = \frac{1}{k} \sum_{i=1}^k \begin{cases} 1 & \text{if } m_i(A) = m_i(B) \\ 0 & \text{otherwise} \end{cases}$$

# hash  
fun

$$\Pr \left[ \left| \hat{JSS}_k(A, B) - \hat{JSS}_k(A, B) \right| \geq \epsilon \right]$$

$$\leq 2 \exp \left( \frac{-2 \epsilon^2}{k \cdot \frac{1}{k^2}} \right)$$

$$k \approx \frac{1}{2 \epsilon^2} \approx 200 \rightarrow 1000 = k$$

### Fast Min-Flashing

$\{0, 1\} \leftarrow \text{rand-bit}()$

$[0, 1] \leftarrow \text{rand-unif}()$

$h: [n] \rightarrow [100 \cdot n \log n]$

for  $x \in S$

$(v_1 = \infty, v_2 = \infty, \dots, v_k = \infty)$

for  $j = 1$  to  $k \log(g)$  ← all entries in vector

if  $(h_j(x) < v_j)$

$v_j \leftarrow h_j(x)$

$v_j \leftarrow x$

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

$(\infty, \infty)$
$(7, 2)$
$(3, 2)$
$(3, 1)$

a  
b  
c

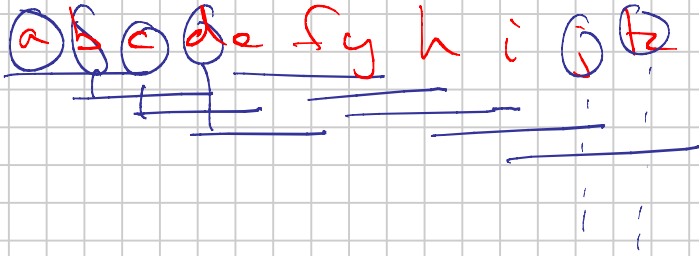
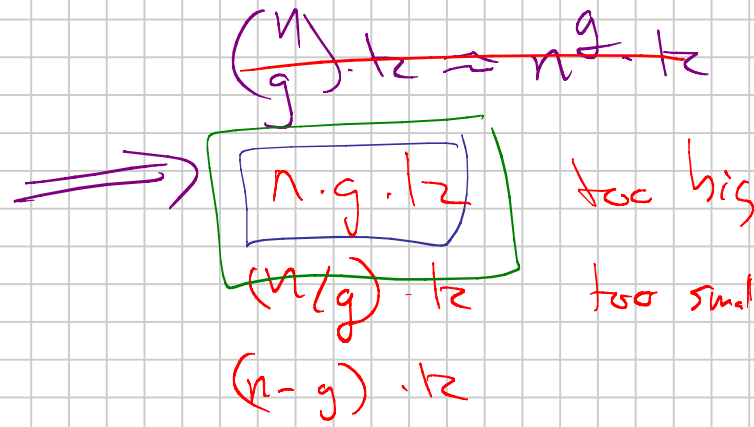
	$h_1$	$h_2$
a	7	2
b	3	10
c	22	1

Runtime  $O(|S| \cdot k)$

Document  
n words

g - word-grams

k - hash fxns



~~$(n - \frac{n}{g}) \cdot k$~~

$O(n \cdot k)$