

LS: Min Hashing

Jan 22
2025

Jeff Phillips

Input: Document \equiv string words, or characters.

\swarrow n-grams
 \rightarrow Sets

$D_1 \rightarrow S_1 \rightarrow V_1$
 $D_2 \rightarrow S_2 \rightarrow V_2$
 $D_3 \rightarrow S_3 \rightarrow V_3$

$V_j \in \mathbb{R}^{1000}$

Jaccard Sim

$$JS(S_1, S_2) = \frac{|S_1 \cap S_2|}{|S_1 \cup S_2|}$$

min hashing \rightarrow

vector

\swarrow
could
have

$D_1 \rightarrow V_1$
 $D_2 \rightarrow V_2$

bag-of-words

$V_j \in \mathbb{R}^{100,000}$

w/ bigram $\mathbb{R}^{(100,000)^2}$

Round-about Min-Hashing

Set \rightarrow matrix \rightarrow vector

Element	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

permutation
2, 6, 4, 3, 1, 5

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

$$S_2 = \{3\}$$

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

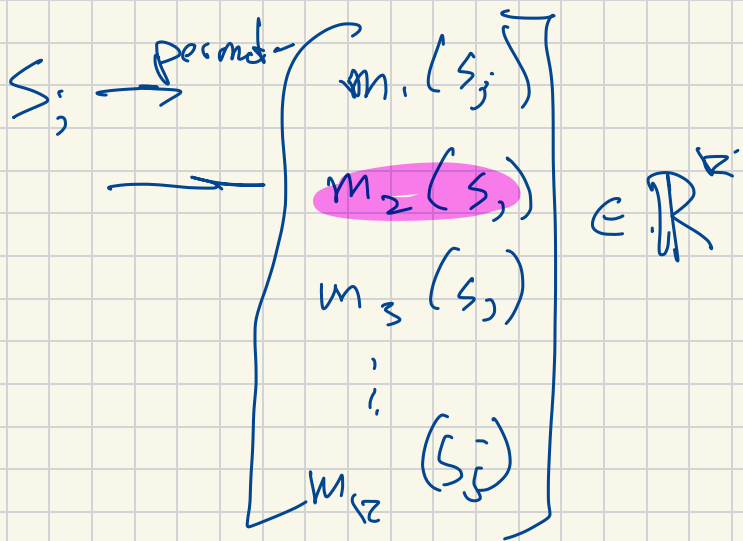
$$S_4 = \{1, 4, 6\}$$

Element	S_1	S_2	S_3	S_4
2	1	0	1	0
6	0	0	1	1
4	0	0	0	1
3	0	1	1	0
1	1	0	0	1
5	1	0	0	0
m_i	2	3	2	6

Map Set to vector \mathbb{R}^k

$k = \# \text{ iterations}$

$\hookrightarrow \# \text{ permutations}$



$m_i(S_j)$

single value representation of

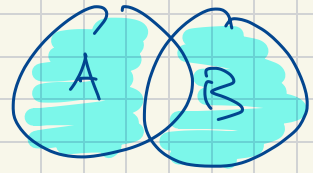
S_j using ith permutation.

Value is first row w/ 1

$$\mathbb{I}_{S_1, S_2} = \frac{1}{k} \sum_{i=1}^k \begin{cases} 1 & \text{if } (v_i[S_1] = v_i[S_2]) \\ 0 & \text{o.w.} \end{cases}$$

$$E[\underbrace{\sum}_\text{perm} \hat{S}_i(s_1, s_2)] = \sum S(s_1, s_2)$$

$$Pr[\underbrace{m}_\text{perm}(s_1) = m(s_2)] = \sum S(s_1, s_2)$$



$$\sum S(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A \cap B| + |A \Delta B|}$$

$A \cap B \Rightarrow \#$
element
in both

which element $[m]$
can be first non-zero $\equiv A \cup B$
in either A or B.

$A \Delta B \Rightarrow \#$
element in
1 or other
not both

which element $[m]$
can give collision $\Rightarrow m_p(A) = m_p(B)$
 $\equiv A \cap B$

Chernoff Hoeffding

k R.V.s X_1, \dots, X_k iid $X_i \in \{0, 1\}$

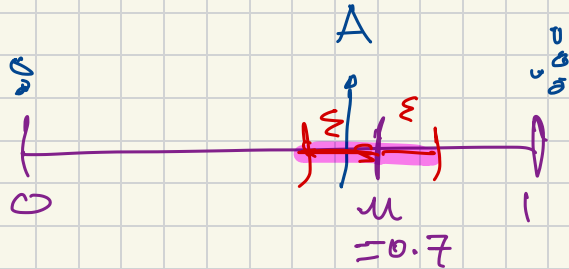
know $E[X_j] = \mu$

$$A = \frac{1}{k} \sum_{j=1}^k X_j$$

$$E[A] = \mu$$

$$\Pr[|A - \mu| > \epsilon] \leq 2 \exp\left(-\frac{\epsilon^2 k}{2}\right) = \delta$$

Probably Approx Correct
PAC



Solve for $k = \frac{2}{\epsilon^2} \log \frac{2}{\delta}$

$$\begin{aligned} \epsilon = 0.05 &\Rightarrow k = 800.8 \\ \delta = 0.01 &\Rightarrow 6400 \end{aligned}$$

Define hash function

$$h_i : [|\Sigma|] \rightarrow [N]$$

\uparrow
things
in all sets

$$\rightarrow [N]$$

random (w/ salt a_1, \dots, a_k , s)
maps uniformly to $[N]$
 $N = |\Sigma|^k$

only need comparison operator

once chosen $h_i \sim \mathcal{H}$, then deterministic.

convert set S to vector $v \in \mathbb{R}^k$

$$h_1, h_2, \dots, h_k \sim \mathcal{H}$$

```
for ( $x_i \in S$ )  
  for ( $j = 1$  to  $k$ )  
    if ( $h_j(x_i) < v_j$ )  
       $v_j \leftarrow h_j(x_i)$ 
```

init $v = (N+1, N+1, \dots, N+1) \in \mathbb{R}^k$

Fast Min Hash

Next

Scenario

$Z = \begin{cases} 1 \text{ billion documents} \\ 1 \text{ billion webpages} \\ 7 \text{ million Wikipedia articles} \end{cases}$

query document

find element $z \in Z$ so $J_S(q, z)$ largest.

quickly w/ scanning all documents

↳ w/ Locality Sensitive Hashing
LSH