# Min Hashing

**Messg Data**
- web page
- news article
- email

$\longrightarrow$ text $\longrightarrow$ vector, bag-of-words $\longrightarrow$ $\mathbb{R}^d$

$\longrightarrow$ Set, b-gram $\longrightarrow$ $S = \{s_1, \ldots s_p\}$ $\xrightarrow[\text{min hashing}]{\text{vector}}$ $[n]^k$

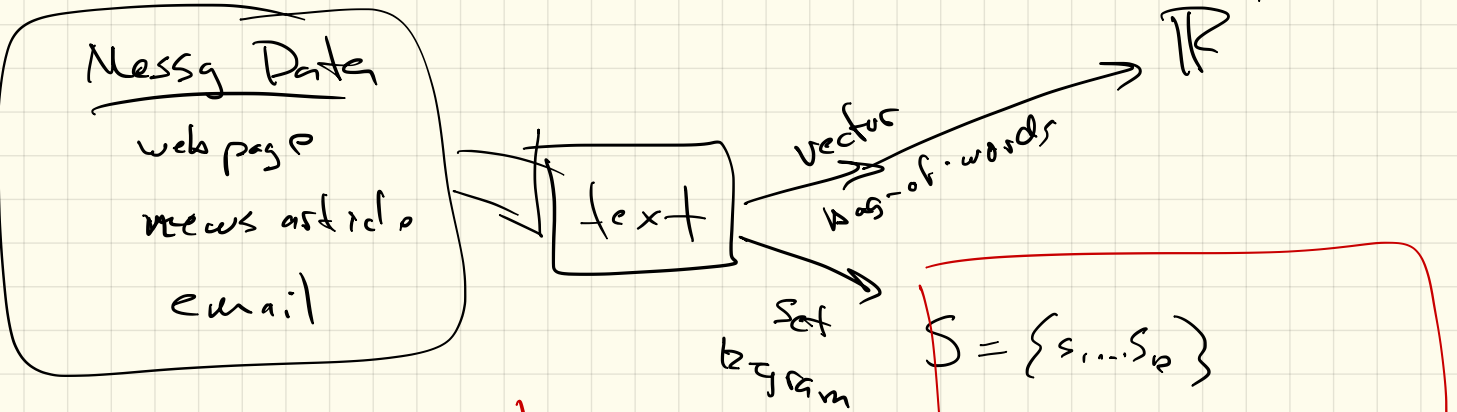**Jaccard Similarity**

$A = \{0, 1, 2, 5\}$

$B = \{2, 3, 5, 6\}$

$$JS(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|\{2, 5\}|}{|\{0, 1, 2, 3, 5, 6\}|} = \frac{2}{6} = \frac{1}{3}$$

Family Hash Functions $\mathcal{H}$

all $h$, where $\sigma$ perm
$[n] \Rightarrow [n]$

Randomly draw $h_{\sigma_i} \in \mathcal{H}$

then $h_{\sigma_i}$ deterministic

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

$S_2 = \{3\}$

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

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

Sets $S_i \subset [n]$

$$h_{\sigma_i} : [n] \to [n]$$
D ← domain

2 ← important!
has order

$h_{\sigma_1}(1) \to 7$

$h_{\sigma_1}(2) \to \boxed{3} \quad 8$

$h_{\sigma_1}(5) \to 4 \quad \boxed{1}$

$h_{\sigma_2}$
2

$$g_1(\hat{S}) = \min_{s \in \hat{S}} h_{\sigma_1}(s)$$

e.g. $g_1(\hat{S}_1) = 3$

$\sigma_i$ = permutation

| domain | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | $[n]=[10]$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $h_{\sigma_1}$ | 7 | 3 | 2 | 6 | 4 | 10 | 1 | 9 | 8 | 5 | → 6 |
| $h_{\sigma_2}$ | 2 | 8 | 6 | 7 | 1 | 5 | 4 | 10 | 9 | 3 | → 2 |

Way to go from

Set $S_i \subset [n]$

$g_1(S_i) \in [n]$

$\times k$

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

$S_2 = \{3\}$

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

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

$$g_1(S_i) \to v_1$$
$$g_2(S_i) \to v_2$$
$$\vdots$$
$$g_k(S_i) \to v_k$$

$$\Rightarrow \quad V = (v_1, v_2, \dots v_k) \in [n]^k$$

$$V(S_1) = (3, 1, , 6 \dots, x)$$

$$V(S_4) = (6, 2, 6 \dots x)$$
$$\qquad\qquad 0 \quad 0 \quad 1$$

$$\text{Apr} \quad \widehat{JS}(S_1, S_4) = \frac{1}{k} \sum_{i=1}^{k} \begin{cases} 1 & \text{if } v_i(S_1) = v_i(S_4) \\ 0 & \text{o.w.} \end{cases}$$

$$= \frac{1}{3}$$

For any two sets $S_1$, $S_2$

$h_{\sigma_1}, h_{\sigma_2}, \ldots h_{\sigma_k} \underset{iid}{\sim} )-($

$$E\left[\hat{JS}(S_1, S_2)\right] = JS(S_1, S_2)$$

---

$$E\left[\hat{JS}(S_1, S_2)\right] = E\left[\frac{1}{k}\sum_{i=1}^{k} \mathbb{1}\left(g_i(S_1) = g_i(S_2)\right)\right]$$

$$= \frac{1}{k}\sum_{i=1}^{k} E\left[\mathbb{1}\left(g_i(S_1) = g_i(S_2)\right)\right]$$

$$\Pr\left[g_i(S_1) = g_i(S_2)\right] = JS(S_1, S_2)$$

Decompose $[n] \rightarrow A, B, C$

A objects hashed to bag $s \in S_1$ and $s \in S_2$

B objects hashed to bag $s \in S_1$ or $s \in S_2$, not both

C objects hashe to bag $x \in S_1 \cup S_2$

$$JS = \frac{|A|}{|A| + |B|}$$

# Fast Min Hash

$$\hat{g}_i \rightleftarrows g_i : \left(\text{set} \subset [n]\right) \longrightarrow [n]$$

choose (random) hash function

$$f_i : [n] \to [m] \qquad m > n$$

$v_i = \infty \quad \forall i$    Domain

for j=1 to l   $\boxed{S = \{x_1, x_2, \dots x_l\}}$

    for i = 1 to k

     if $\left( f_i(x_j) < v_i \right)$

       $v_i \leftarrow h_i(x_j)$

Return $V = (v_1, v_2, \dots v_k)$

Domain
$$\Omega = [n]$$

every set
$$S \subset \Omega$$

$$x \in S$$

$$S \in 2^\Omega$$

$\subseteq H$

## How large should $k$ be?

$X_1, X_2, \ldots X_k \underset{iid}{\sim} \mu \qquad X_i \in [0,1]$

$$E[A] = E[X_i] = \mu$$

$$A = \frac{1}{k} \sum_{i=1}^{k} X_i$$

$$\Pr\left[ |A - \mu| > \varepsilon \right] \leq 2 \exp\left(-2\varepsilon^2 k\right) \underset{0.01}{\leq \delta}$$

$$\varepsilon = 0.05$$

$$\delta = 2\exp\left(-2(0.05)^2 k\right)$$

$$\ln\left(\frac{\delta}{2}\right) = -2\left(\frac{1}{10}\right)^2 k$$

$$\ln\left(\frac{2}{\delta}\right) = 2\left(\frac{1}{10}\right)^2 k \implies k = \frac{400}{2} \ln(200)$$

$$k = 200 \ln(200) = 1060$$

$h_{\sigma_1}(1) \to 7$

$h_{\sigma_1}(2) \to \boxed{3}$   $h_{\sigma_2}$   $2$   $8$   $\boxed{1}$

$h_{\sigma_1}(5) \to 4$

$g_1(s) = \min_{s \in \hat{s}} h_{\sigma_1}(s)$

e.g. $g_1(s_1) = 3$

$\sigma_i = \text{permutation}$

domain $1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9\ 10$   $[n] = [10]$

$h_{\sigma_1}$  $7\ 3\ 2\ 6\ 4\ 10\ 1\ 9\ 8\ 5 \to 6$

$h_{\sigma_2}$  $2\ 8\ 6\ 7\ 1\ 4\ 10\ 9\ 3 \to 2$

$h_3$  $1\ 6\ 4\ 2\ 7\ 8\ 10\ 9\ 5\ 3$

A B C B B B C C C C

|  | $h_1$ | $h_2$ | $h_3$ |
|---|---|---|---|
| $S_1$ | $7\ 3\ 4$ | $2\ 8\ 1$ | $1\ 6\ 7$ |
| $S_2$ | $7\ 6\ 10$ | $2\ 7\ 5$ | $1\ 2\ 8$ |
|  | 0 | 8 | 1 |