# Min - Hashing

## Jaccard Sim

Sets $\longrightarrow$ Matrix $\longrightarrow$ Binary Vector

hashing

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

$A = \{0, 1, 2, 5, 6\}$

$B = \{0, 2, 3, 5, 7\}$

$$= \frac{|\{0, 2, 5\}|}{|\{0,1,2,3,5,6,7\}|} = \frac{3}{7}$$

Compare

$\frac{3}{7}$ LSH

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

$$S_2 = \{3\}$$

$$S_3 = \{2, 3, 4, 5\}$$

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

$$JS(S_1, S_3) = \frac{|\{2, 5\}|}{|\{1, 2, 3, 4, 5\}|} = \frac{2}{5}$$

$$\boxed{E\left[\hat{JS}(S_i, S_j)\right] = JS(S_i, S_j)}$$

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

mostly 0s

$\xrightarrow{\quad\times k \text{ times}\quad}$ random reorder rows

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

$$m \quad 2 \quad 3 \quad 2 \quad 6$$

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

$k$ random reorders $\quad j = 1, 2, \ldots k$

Set $S_{i=3}$ value $\left( m_1(S_i), m_2(S_i), \ldots m_p(S_i) \right)$

$$V_i = V_3 = \left( 2, \ \textcircled{7}, \ 10023, \ \ldots, \ \textcircled{18} \right)$$

$$S_{i'=1} \implies V_{i'} = V_1 = \left( 3, \ \textcircled{7}, \ 42, \ \ldots, \ \textcircled{18} \right)$$

$$\phantom{S_{i'=1} \implies V_{i'} = V_1 = \left( } \ 0 \quad\ \ 1 \quad\ \ \ 0 \quad \ldots \quad 1$$

$$\hat{JS}_j(S_{i'}, S_i) = \begin{cases} 1 & \text{if } \ m_j(S_{i'}) = m_j(S_i) \\ 0 & \text{o.w.} \end{cases}$$

$$\overline{JS} = \frac{1}{k} \sum_{j=1}^{k} \hat{JS}_j(S_{i'}, S_i)$$

## Why $E[\hat{JS}(S_1, S_2)] = JS(S_1, S_2)$

3 types of rows

$(T_x) = \#$ rows $S_1$ and $S_2$ have 1

$(T_y) = \#$ rows exactly 1 of $S_1, S_2$ have 1
$\qquad\qquad\qquad\qquad\qquad$ other has 0

$(T_z) = \#$ rows $S_1$ and $S_2$ have 0.

$$JS(S_1, S_2) = \frac{T_x}{T_x + T_y} \quad\Longrightarrow\quad \text{Can ignore } (T_z)$$

Collision $m(S_1) = m(S_2)$ iff of $T_x$ and $T_y$ rows
$\qquad\qquad\qquad$ a $T_x$ is at the top.

$\text{Prob}[\text{collision}] = E[\hat{JS}] = \frac{T_x}{T_x + T_y} = JS$

# Top k - Sketches

$k = 2$

$S_i$

| | |
|---|---|
| 1 | 0 |
| 5 | 1 |
| 2 | 0 |
| 3 | 1 |
| 6 | 1 |
| 4 | 0 |

$$m(S_i) = (5, 3)$$

# Fast Min Hashing Algorithm

Replace re-orders w/ hash fxns.

Hash fxns $h_1, h_2, \ldots, h_k$

$$h_j : \{\text{element of set}\} \longrightarrow [n]$$

$$= \{1, 2, \ldots, 1024\}$$

$S \rightarrow V = (v_1, v_2, \ldots, v_k)$

$v_2$ replace $m_2(S)$
$v_j$ replace $m_j(S)$

```
for i ∈ S
  for j = 1 to k    (all hash)
    if (h_j(i) < v_j)
      v_j ← h_j(i)
```

$$S = \{1, 3, 6\}$$

$$
\begin{array}{c|c}
 & S \\
\hline
1 & 1 \\
2 & 0 \\
3 & 1 \\
4 & 0 \\
5 & 0 \\
6 & 1 \\
\end{array}
$$

$$
h_1 \begin{vmatrix}
1 \to 2 \\
2 \to 4 \\
3 \to 6 \\
4 \to 4 \\
5 \to 1 \\
6 \to 3 \\
\end{vmatrix}
$$

$$
h_2 \begin{vmatrix}
1 \to 3 \\
2 \to 5 \\
3 \to 1 \\
4 \to 6 \\
5 \to 2 \\
6 \to 6 \\
\end{vmatrix}
$$

$$V_1 = \infty$$

$1 \quad V_1 = 2$

$3 \quad V_1 = 2 \qquad 2 < 6$

$6 \quad V_1 = 2 \qquad 2 < 3$

$$V_2 = 1$$

hash fxns are fixed for all $s \in S$

$$H : \{h_1, \dots h_k\}$$

$$f_H(S) \to \mathbb{R}^k \equiv V_S$$

# Central Limit Theorem

$$X = \{x_1, x_2, \dots x_r\} \qquad r \text{ random variables}$$
$$\text{iid.}$$

$$A = \frac{1}{r} \sum_{i=1}^{r} x_i$$

$$E[A] = E[x_i]$$

$$Var(A) = \frac{Var[x_i]}{r}$$



$\sqrt{Var(A)}$

$\delta/2$       $\epsilon$   $E[x]$   $\delta/2$

Value A

Converge to Normal

# Probably Aprox Correct

$$\Pr\left[\,|A - E[X_i]| \geq \varepsilon\,\right] < \delta$$

estimate

what I want

error tolerance

probabilty of failure

Want $\varepsilon, \delta$ small

Algorithm has $r$ steps

# Chernoff - Hueffding Bound

iid $X_1, \ldots X_r$        $A = \frac{1}{r} \sum_{i=1}^{r} X_i$

$r = k$

$\hat{JS}$    $\Delta = b - a$    s.t. $X_i \in [a, b]$

$\Delta = 1$                                       $\underset{0}{}$  $\underset{1}{}$

$$Pr\left[ | A - E[X_i] | > \varepsilon \right] \leq 2 \exp\left( \frac{-2 r \varepsilon^2}{\Delta^2} \right)$$

$$Pr\left[ | \bar{JS} - JS | > \underset{\varepsilon}{0.1} \right] \leq 2 \exp\left( \frac{-2 k (0.1)^2}{1} \right)$$

$k = 500$

$$\leq 2 \exp\left( \frac{-2k}{100} \right) = 2\exp\left( \frac{-1000}{100} \right)$$

$$= 2 e^{-10}$$