

BDA  
Assignment-1

group- C7  
(Prince, Bhautik, Lyon)

1. Shingles:-

(a) Here 2-shingles for ABRACADABRA are;

$$\mathcal{D}_1 = \{AB, AC, AD, BR, CA, DA, RA\}$$

$$\text{so } |\mathcal{D}_1| = 7$$

(b) 2-shingles for BRICABRAC are;

$$\mathcal{D}_2 = \{AB, AC, BR, CA, IC, RA, RI\}$$

$$\text{so } |\mathcal{D}_2| = 7$$

(c) The common 2-shingles are;

$$[\mathcal{D}_1 \cap \mathcal{D}_2] = \{AC, AB, BR, CA, RA\}$$

$$\therefore |\mathcal{D}_1 \cap \mathcal{D}_2| = 5$$

(d) The Jaccard similarity is defined as;

$$\text{SIM}(\mathcal{D}_1, \mathcal{D}_2)$$

$$\text{SIM}(\mathcal{D}_1, \mathcal{D}_2) = \frac{|\mathcal{D}_1 \cap \mathcal{D}_2|}{|\mathcal{D}_1 \cup \mathcal{D}_2|} = \frac{5}{9}$$

$$\boxed{\text{SIM}(\mathcal{D}_1, \mathcal{D}_2) = \frac{5}{9}}$$

## 2. Minhashing:-

(a) Permuting rows of characteristic matrix  $X$

① using permutation vector,  
 $(4, 5, 0, 2, 3, 1)^T$

A characteristic matrix of four sets  $(C_1, C_2, C_3, C_4)$  over universal set  $\{0, 1, 2, 3, 4, 5\}$  and permutation of its rows,

$0 \rightarrow 4, 1 \rightarrow 5, 2 \rightarrow 0, 3 \rightarrow 2, 4 \rightarrow 3, 5 \rightarrow 1$   
is given as:

Row	$C_1$	$C_2$	$C_3$	$C_4$		Row	$C_1$	$C_2$	$C_3$	$C_4$
0	0	1	1	0		4	1	0	1	0
1	1	0	1	1		5	0	1	0	0
2	0	1	0	1	$\rightarrow$	0	0	1	1	0
3	0	0	1	0		2	0	1	0	1
4	1	0	1	0		3	0	0	1	0
5	0	1	0	0		1	1	0	1	1

The minhash function  $h_\pi$  on  $C$  is defined by,

$$h_\pi(c) = \min_{i \in \{1, \dots, m\}} \{ \pi(i) \mid c[i] = 1 \}$$



Where we are given;

- a characteristic matrix with  $m$  rows = 6 and column  $c$ .
- a permutation  $\pi$  on the rows, that is  $\pi: \{1, \dots, m\} \rightarrow \{1, \dots, m\}$  is a bijection. (1 corresponds to 0)

$$\pi: 4 \rightarrow 0, 5 \rightarrow 1, 0 \rightarrow 2, 2 \rightarrow 3, 3 \rightarrow 4, 1 \rightarrow 5$$

$$\text{so; } h_{\pi}(c_1) = 0, h_{\pi}(c_2) = 1, h_{\pi}(c_3) = 0, h_{\pi}(c_4) = 3$$

(2) Using the Permutation vectors;  $(3, 1, 0, 5, 2, 4)^T$

$$\pi: 3 \rightarrow 0, 1 \rightarrow 1, 0 \rightarrow 2, 5 \rightarrow 3, 2 \rightarrow 4, 4 \rightarrow 5$$

Permutation of it's row is given by:

Row	$c_1$	$c_2$	$c_3$	$c_4$
0	0	1	1	0
1	1	0	1	1
2	0	1	0	1
3	0	0	1	0
4	1	0	1	0
5	0	1	0	0

→

Row	$c_1$	$c_2$	$c_3$	$c_4$
3	0	0	1	0
1	1	0	1	1
0	0	1	1	0
5	0	1	0	0
2	0	1	0	1
4	1	0	1	0

$$\text{so; } h_{\pi}(c_1) = 1, h_{\pi}(c_2) = 2, h_{\pi}(c_3) = 0, h_{\pi}(c_4) = 1.$$

(counting from 0).

(b)

Row	$s_1$	$s_2$	$s_3$	$s_4$	$h_1(x)$ $2x+1 \pmod 7$	$h_2(x)$ $3x+2 \pmod 7$	$h_3(x)$ $5x+2 \pmod 7$
0	0	1	0	1	1	2	2
1	0	0	0	1	3	5	0
2	1	1	0	0	5	1	5
3	0	0	1	0	0	4	3
4	0	1	1	0	2	0	1
5	1	0	0	0	4	3	6
6	1	0	1	0	6	6	4

First iteration gives:

$$SIG_{12} = SIG_{14} = \min(\infty, h_1(0)) = 1$$

$$SIG_{22} = SIG_{24} = \min\{\infty, h_2(0)\} = 2$$

$$SIG_{32} = SIG_{34} = \min(\infty, h_3(0)) = 2$$

	$s_1$	$s_2$	$s_3$	$s_4$
$h_1$	$\infty$	1	$\infty$	1
$h_2$	$\infty$	2	$\infty$	2
$h_3$	$\infty$	2	$\infty$	2

Similarly second iteration:

	$s_1$	$s_2$	$s_3$	$s_4$
$h_1$	$\infty$	1	$\infty$	1
$h_2$	$\infty$	2	$\infty$	2
$h_3$	$\infty$	2	$\infty$	0



Third iteration;

	$s_1$	$s_2$	$s_3$	$s_4$
$h_1$	5	1	$\infty$	1
$h_2$	1	1	$\infty$	2
$h_3$	5	2	$\infty$	0

fourth iteration

	$s_1$	$s_2$	$s_3$	$s_4$
$h_1$	5	1	0	1
$h_2$	1	1	4	2
$h_3$	5	2	3	0

~~fourth~~ iteration

	$s_1$	$s_2$	$s_3$	$s_4$
$h_1$	5	1	0	1
$h_2$	1	0	0	2
$h_3$	5	1	1	0

~~fifth~~ sixth iteration

	$s_1$	$s_2$	$s_3$	$s_4$
$h_1$	4	1	0	1
$h_2$	1	0	0	2
$h_3$	5	1	1	0

Seventh / final iteration

	$s_1$	$s_2$	$s_3$	$s_4$
$h_1$	4	1	0	1
$h_2$	1	0	0	2
$h_3$	4	1	1	0

(b) here  $h_1$  and  $h_3$  are permutations

(c) Jaccard similarities;

$$\begin{aligned} \text{SIM}(s_1, s_2) &= \frac{|s_1 \cap s_2|}{|s_1 \cup s_2|} = 0 = \text{SIM}(s_1, s_3) \\ &= \text{SIM}(s_1, s_4) \\ &= \text{SIM}(s_3, s_4) \end{aligned}$$

$$\text{SIM}(s_2, s_3) = \frac{2}{3},$$

$$\text{SIM}(s_2, s_4) = \frac{1}{3}.$$

True Jaccard Similarity:

$$\text{SIM}(s_1, s_2) = \frac{x}{x+y} = \frac{1}{4+1} = \frac{1}{5}.$$

$$\text{SIM}(s_1, s_3) = \frac{1}{5}, \quad \text{SIM}(s_2, s_4) = \frac{1}{4}$$

$$\text{SIM}(s_1, s_4) = 0 \quad \text{SIM}(s_3, s_4) = 0$$

$$\text{SIM}(s_2, s_3) = \frac{1}{5}$$

for pair  $(s_1, s_4) \cap (s_3, s_4)$  it's equal

Pair  $(s_2, s_3)$  has a difference of  $\frac{1}{12}$

which is small, so this pair has close estimated and True Jaccard Similarity.

(d) Benefits of using hash functions instead of permutations;

- Requires less memory.

- easy to generate

- No memorization required.



- Permutation require iteration over whole Problem space while for hash function we can apply iteratively
- Some hash functions come with the guarantee of uniformity, randomness
- Permutation is time consuming for large number of rows while hash function is easier as it is not a bijection.