

## Assignment - 2

### Question 1

	$C_1$	$C_2$	$C_3$	$C_4$	$C_5$	$C_6$	$C_7$
Band 1	1	2	1	1	2	5	4
	2	3	4	2	3	2	2
Band 2	3	1	2	3	1	3	2
	4	1	3	1	2	4	4
Band 3	5	2	5	1	1	5	1
	6	1	6	4	1	1	4

$$\text{Band 1} = (C_1, C_4), (C_2, C_5)$$

$$\text{Band 2} = (C_1, C_6)$$

$$\text{Band 3} = (C_1, C_3), (C_4, C_7)$$

### Question 2

$$M = 50\% \Rightarrow \text{false positive}$$

$$N = 20\% \Rightarrow \text{false negative}$$

$R$	$b$	FP(M)	FN(N)
1	24	0.9952	0.0000005
2	12	0.3872	0.0316
3	8	0.0622	0.343
4	6	0.0095	0.6789
6	4	0.002	0.9389
8	3	$7.67E^{-6}$	0.98832
12	2	$8.192E^{-9}$	0.9995
24	1	0	0.99999

Y: rows  
b: bins

$$\left[ 1 - (1 - S^b) \right]^b$$

for FP

for FN  
without the  
2

a) 1:1  $\rightarrow$  M:N  $\rightarrow$

- 0.995277
- 0.418966
- 0.405845
- 0.688495
- 0.988334
- 0.999511
- 0.99999

~~b) 1:10  $\rightarrow$  M:10N  $\rightarrow$~~

1:10 M:10N	1:100 M:100N	1:1000 M:1000N
0.99527	0.99528	0.99533
0.70405	0.55492	32.06364
3.49832	32.4231	343.6711
6.79890	67.9029	678.9473
9.38975	93.8952	938.9499
9.88327	98.8327	988.3269
9.995117	99.9511	999.511
9.99999	99.9999	999.9997

Question 3

Find the set of 2 shingles for documents given

document 1 : ABRACADABRA

Set (A)  $\Rightarrow$  Shingle of size 2

$\hookrightarrow \{AB, BR, RA, AC, CA, AD, DA, AB, BR, RA\}$

So  $\Rightarrow \{AB, BR, RA, AC, CA, AD, DA\}$

document 2 : BRICABRAC

Set B  $\Rightarrow \{BR, RI, CA, AB, BR, RA, AC\}$

~~Union~~

① 2 shingles document 1 has  $\Rightarrow 7$

② 2 shingles document 2 has  $\Rightarrow 7$

③ 2 shingles they have in common

doc-1  $\Rightarrow AB, BR, RA, AC, CA, AD, DA$

doc-2  $\Rightarrow BR, RI, CA, AB, BR, RA, AC$

so 5 shingles

Jaccard similarity  $\Rightarrow \frac{\text{Intersection}}{\text{Union}}$

$$= \frac{5}{7+7-5} = \frac{5}{9}$$



Question 4:

Compute the Jaccard similarity between each pair of columns.

	C1	C2	C3	C4
R1	0	1	1	0
R2	1	0	1	1
R3	0	1	0	1
R4	0	0	1	0
R5	1	0	1	0
R6	0	1	0	0

$$J(C1, C2) = 0/5 = 0$$

$$J(C1, C3) = 2/4 = 0.5$$

$$J(C1, C4) = 1/3 = 0.33$$

$$J(C2, C3) = 1/6 = 0.167$$

$$J(C2, C4) = 1/4 = 0.25$$

$$J(C3, C4) = 1/5 = 0.2$$

Question 5:

	C1	C2	C3	C4
R4	0	0	1	0
R6	0	1	0	0
R1	0	1	1	0
R3	0	1	0	1
R5	1	0	1	0
R2	1	0	1	1

Minhash value of C1 = R5

Minhash value of C2 = R6

Minhash value of C3 = R4

Minhash value of C4 = R3