

# Homework Assignment 2

---

## Problem 1

$$E(sim(S, T)) = \begin{cases} 1 & \text{if } 0 < n \leq m; \\ \sum_{k=m}^n \frac{2m-k}{k} \frac{\binom{m}{2m-k} \binom{n-m}{k-m}}{\binom{n}{m}} & \text{if } m < n \leq 2m-2; \\ \sum_{k=m}^{2m-1} \frac{2m-k}{k} \frac{\binom{m}{2m-k} \binom{n-m}{k-m}}{\binom{n}{m}} & \text{if } n > 2m-2; \end{cases}$$


---

## Problem 2

### Question 1

```
{
'Even hash them ',
'hash them four bytes each, ',
'four bytes each, space needed ',
'space needed store ',
'store still roughly four times ',
'still roughly four times space taken ',
'space taken document ',
'document '
}
```

### Question 2

n-k+1

---

## Problem 3

Map phase:

Map(index\_of\_band, index\_of\_doc, shingles\_of\_every\_doc) -> <index\_of\_band, index\_of\_doc,

hash\_per\_band>

Inside every Mapper hash table for every document is calculated within

same band. Meanwhile indexes of documents are only be preserved.

Reduce phase:

Reduce(index\_of\_doc, hash\_table\_per\_band) -> <index\_of\_doc, hash\_table\_per\_doc>

For every Reducer hash tables per band are grouped and ordered according to document id

---

#### Problem 4

$$\begin{aligned}
 & 1 - (1 - s^r)^b \\
 &= 1 - (1 - s^{\frac{r}{2}})^b (1 + s^{\frac{r}{2}})^b \\
 &= 1 - \lim_{n \rightarrow \infty} \prod_{i=1}^n e^{\frac{s^r b}{2^i}} \\
 &= 1 - e^{s^r b \lim_{n \rightarrow \infty} \sum_{i=1}^n \frac{1}{2^i}} \\
 &= 1 - e^{s^r b}
 \end{aligned}$$


---

#### Problem 5

1. precise threshold : 0.569

estimate threshold : 0.607

relative difference(|pt-et|/pt) : 0.066

2. precise threshold : 0.406

estimate threshold : 0.464

relative difference : 0.143

3. precise threshold : 0.880

estimate threshold : 0.890

relative difference : 0.011

The estimate threshold approaches to precise threshold when the value of formula is exactly 1/2, especially when the value of b as well as r are sufficiently great.

## Problem 6

1.  $\max(x, y)$  is not a distance measure. When  $x \neq 0$ ,  $\max(x, x) = x \neq 0$  Identity not holds.

2.  $\text{diff}(x, y)$  is a distance measure. Because

1. Non-negativity:  $\text{diff}(x, y) = |x - y| \geq 0$ , if and only if  $x = y$  the equality holds.

2. Symmetricity:  $\text{diff}(x, y) = |x - y| = |y - x| = \text{diff}(y, x)$

1. Identity:  $\text{diff}(x, y) = |x - y| = 0$ , iff  $x = y$

2. Triangle inequality:  $\text{diff}(x, y) + \text{diff}(y, z) = |x - y| + |y - z|$  case 1:  $x = z$   $|x - y| + |y - z| = |x - y| + |y - x| \geq 0 = |x - x| = |x - y| = \text{diff}(x, y)$  case 2:  $x \neq z$ , let  $x < z$  inequality  $\{x - y + z - y = z - x + 2x - 2y > z - x = \text{diff}(x, z)$ , if  $0 < y < x < z$   $\{y - x + z - y = z - x = \text{diff}(x, z)$ , if  $x \leq y < z$   $\{y - z + y - x = z - x - 2z + 2y > z - x = \text{diff}(x, z)$ , if  $x < z \leq y$  so, it's a metric

3.  $\text{sum}(x, y) = x + y$  is not a metric. if  $x \neq 0$   $\text{sum}(x, x) = 2x \neq 0$  Identity not holds.

4. Jaccard distance is a metric. Given three sets A, B and C, Jaccard distance of two sets (e.g. A and B) is defined as  $J(A, B) = 1 - |A \cap B| / |A \cup B|$

1. Non-negativity:  $1 - |A \cap B| / |A \cup B| \geq 0$ , iff  $A = B$  the equality holds.

2. Symmetricity:  $J(A, B) = 1 - |A \cap B| / |A \cup B| = 1 - |B \cap A| / |A \cup B| = J(B, A)$

3. Identity:  $J(A, A) = 1 - |A \cap A| / |A \cup A| = 1 - |A| / |A| = 0$

4. Triangle inequality: see solution below:

$$\begin{aligned}
 & a + A \setminus B \setminus C, b + B \setminus C \setminus A, c + C \setminus A \setminus B \\
 & \alpha + A \cap C \setminus B, \beta + a \cap B \cap C, \gamma + A \cap B \setminus C, \delta + B \cap C \setminus A \\
 & J(A, B) + \frac{a + b + \alpha + \delta}{\alpha + \beta + \gamma + \delta + a + b} \\
 & J(B, C) + \frac{b + c + \alpha + \gamma}{\alpha + \beta + \gamma + \delta + b + c} \\
 & J(A, C) + \frac{a + c + \gamma + \delta}{\alpha + \beta + \gamma + \delta + a + c} \\
 & J(A, B) + J(B, C) \geq J(A, C) \\
 & \Leftrightarrow \frac{a + b + \alpha + \delta}{\alpha + \beta + \gamma + \delta + a + b} + \frac{b + c + \alpha + \gamma}{\alpha + \beta + \gamma + \delta + b + c} \geq \frac{a + c + \gamma + \delta}{\alpha + \beta + \gamma + \delta + a + c}
 \end{aligned}$$

$$\begin{aligned}
&\Leftrightarrow \frac{a+\delta}{\alpha+\beta+\gamma+\delta+a+b} + \frac{c+\gamma}{\alpha+\beta+\gamma+\delta+b+c} + \frac{b+\alpha}{\alpha+\beta+\gamma+\delta+a+b} + \\
&\frac{b+\alpha}{\alpha+\beta+\gamma+\delta+b+c} \geq \frac{a+c+\gamma+\delta}{\alpha+\beta+\gamma+\delta+a+c} \\
&\Leftrightarrow \frac{a+c+\gamma+\delta}{\alpha+\beta+\gamma+\delta+a+b+c} + \frac{b+\alpha}{\alpha+\beta+\gamma+\delta+a+b} + \frac{b+\alpha}{\alpha+\beta+\gamma+\delta+b+c} \geq \\
&\frac{a+c+\gamma+\delta}{\alpha+\beta+\gamma+\delta+a+b+c} \\
&\Leftrightarrow \frac{b+\alpha}{\alpha+\beta+\gamma+\delta+a+b} + \frac{b+\alpha}{\alpha+\beta+\gamma+\delta+b+c} \geq \frac{b}{\alpha+\beta+\gamma+\delta+a+b+c} \\
&\Leftrightarrow \frac{b+\alpha}{\alpha+\beta+\gamma+\delta+a+b+c} + \frac{b+\alpha}{\alpha+\beta+\gamma+\delta+b+c} \geq \frac{b}{\alpha+\beta+\gamma+\delta+a+b+c} \\
&\Leftrightarrow \frac{\alpha}{\alpha+\beta+\gamma+\delta+a+b+c} + \frac{b+\alpha}{\alpha+\beta+\gamma+\delta+b+c} \geq 0
\end{aligned}$$

5. Shortest path is a distance measure Given a graph  $G=(V,E,W)$ , where by  $V$  is the set of vertex,  $E$  is the set of edge which is a sub set of  $V \times V$ ,  $W$  is the set of weight, for each edge there is a non-negative weight.

Define  $\text{len}(V') = \sum(w' \in W)$  where  $V' \subseteq V$ ,  $w$ 's are with  $V'$  corresponded weight.

Define  $\text{sp}(x, y): \forall x \in V, \forall y \in V: (V(x, y) \vee (\exists z \in V: \text{sp}(x, z) \wedge V(z, y)) \wedge (\neg \exists z \in V: \text{len}(\text{sp}(x, z)) + \text{len}(\text{sp}(z, y)) < \text{len}(\text{sp}(x, y))))$ .

1. Non-negativity: Since length of shortest path is defined as sum of weight of subset of edge, and weight edges is non-negative, length of shortest path is non-negative.
2. Symmetricity: Since the graph  $G$  is a undirected graph,  $\text{length}(\text{sp}(x, y)) = \text{length}(\text{sp}(y, x))$
3. Identity: Obviously  $\text{len}(\text{sp}(x, x)) = 0 \quad \forall x \in V$
4. Triangle inequality:

given  $\text{sp}(x, y) (\forall x \in V, \forall y \in V)$  is a shortest path between  $x$  and  $y$ .

assume  $\exists z \in V$ , let  $\text{len}(\text{sp}(x, z)) + \text{len}(\text{sp}(z, y)) < \text{len}(\text{sp}(x, y))$ .

Then  $z$  either on  $\text{sp}(x, y)$ , which meet a contradiction of the definition of length of a path, or not on  $\text{sp}(x, y)$ , which leads to another contradiction that  $x \rightarrow z \rightarrow y$  is the shortest path, instead of  $x \rightarrow y$ .

$\therefore \text{len}(\text{sp}(x, z)) + \text{len}(\text{sp}(z, y)) \geq \text{len}(\text{sp}(x, y))$