

### Question 1 :-

① Inverted index posting lists for terms 'Cat' and 'dog' as follows :-

Cat : 234569:1, 234578:1, 234839:1

dog : 234569:1, 234578:1, 234879:1

② Here is the compressed form for 'dog' :-

dog : 234569:1, 9:1, 301:1 [using gaps]

In binary :-

$$234569 = b111001\ 0100\ 0100\ 1001$$

$$1 = b1$$

$$9 = b1001$$

$$1 = b1$$

$$301 = b1\ 0010\ 1101$$

$$1 = b1$$

VB encoding :-

00001110 00101000 11001001

10000001

10001001

10000001

00000010 10101101

10000001

③ We compute the document vector magnitude as follows

Document ID - 234569

Text : phrase fight cat dog reflect natural  
tendency relationship two species  
antagonistic two species friend

maxf : 2 (two, species)

$$\therefore \text{vector magn.} : \sqrt{\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2}$$
$$= 2.12$$

Doc ID: 234578

Text : Dog cat bad relationship

maxf : 1

$$\text{vec. magn.} : \sqrt{4} = 2$$

DOCID : 234839

Text : Cat fury

maxf : 1

$$\text{vec. magn.} : \sqrt{2} = 1.41$$

DOC ID : 234879

Text : Dog man best friend

maxf : 1

$$\text{vec. magn.} : \sqrt{4} = 2$$

q : cat dog

maxf : 1

$$\text{vec. magn.} : \sqrt{2} = 1.41$$

$$\therefore \cos(q, d_{234569}) = \frac{\left(\frac{1}{2} \times \frac{1}{1}\right) + \left(\frac{1}{2} \times \frac{1}{1}\right)}{1.41 \times 2.12}$$
$$= 0.33$$

$$\cos(q, d_{234578}) = \frac{1+1}{1.41 \times 2}$$
$$= 0.71$$

$$\cos(q, d_{234839}) = \frac{1}{1.41 \times 1.41}$$
$$= 0.50$$

$$\cos(q, d_{234879}) = \frac{1}{1.41 \times 2}$$
$$= 0.35$$

∴ Document 234578 is most similar.

### Question 2:

① The largest integer?

Map(id, i)  
emit(1, i)

Combine(j, list)  
emit(1, max(list))

Reduce(j, list)  
emit(max(list), null)

⑤ The average of all integers ?

Map( $i_d, i$ )  
emit(1,  $i$ )

Combine( $j, \text{list}$ )  $\rightarrow n = \text{sizeof}(\text{list})$   
emit(1,  $(\sum_n i, n)$ )

Reduce( $j, \text{list}$ )  $\rightarrow [(\frac{\sum_i i}{n_1}, n_1), (\frac{\sum_i i}{n_2}, n_2), \dots]$   
emit( $\frac{\sum_m (\text{sum}_i \times n_i)}{\sum_m n_i}, \text{null}$ )  
 $m = \text{size of}(\text{list})$

⑥ Same set of integers (no duplicate)

Map( $i_d, i$ )  
emit( $i, \text{null}$ )

Reduce( $j, \text{list}$ )  $\rightarrow$  [taking advantage of MapReduce ability to group keys]  
emit( $j$ ,  $\text{null}$ )

⑦ output from ⑥ + ↗

Map2( $i_d, i$ )  
emit(1,  $i$ )

Reduce2( $j, \text{list}$ )  $\rightarrow$  [distinct values]  
emit( $\text{sizeof}(\text{list}), \text{null}$ )

### Question 3:

#### ① Bag Union :

$t$  is tuple,  $tid$  is tuple id,  $x$  is bit indicating two Relation ( $R$  or  $S$ )

Map( $tid, (t, x)$ )  
emit( $t, x$ )

Reduce( $t, li$ )

$ln = [x \text{ for } x \text{ in } li \text{ if } x == R]$   
 $ls = [x \text{ for } x \text{ in } li \text{ if } x == S]$   
 emit( $t, \text{sum}(\text{len}(ln), \text{len}(ls))$ )

defined in question

~~defn in practice test~~ ~~(len(ln) + len(ls))~~ ~~min(len(ln), len(ls))~~

#### ② Bag Intersection:

Map( $tid, (t, x)$ )  
emit( $t, x$ )

Reduce( $t, li$ )

$ln = [x \text{ for } x \text{ in } li \text{ if } x == R]$   
 $ls = [x \text{ for } x \text{ in } li \text{ if } x == S]$   
 $c = \min(\text{len}(ln), \text{len}(ls))$   
 $\text{if}(c > 0)$   
 emit( $t, c$ )

## ② Bag difference

Map ( $t, id, (t, x)$ )

emit ( $t, x$ )

Reduce ( $t, li$ )  $\rightarrow [R, R, S, S, R, S \dots]$

$l_R = [x \text{ for } x \text{ in } li \text{ if } x == R]$

$l_S = [x \text{ for } x \text{ in } li \text{ if } x == S]$

$c = \text{len}(l_R) - \text{len}(l_S)$

if ( $c > 0$ )

emit ( $t, c$ )