# QUIZ-4 SOLUTIONS

| BUCKETS | HASH FUNC 1 $(i+j \bmod 3)$ ① | HASH FUNC 2 $(j \bmod 3)$ ② | Counts [1.5] | | [1.5] | Bitmap | |
|---|---|---|---|---|---|---|---|
| | | | ① | ② | ① | ② | |
| 0 | (1,2) (1,5) (2,4) (4,5) | (1,3) (1,3) (2,3) (2,3) | 4 | 4 | 1 | 1 | |
| 1 | (1,3) (3,4) (1,3) (3,4) | (2,4) (3,4) (3,4) | 4 | 3 | 1 | 1 | |
| 2 | (2,3) (3,5) (2,3) (3,5) | (1,2) (3,5) (3,5) (4,5) (1,5) | 4 | 5 | 1 | 1 | |

**Multihash**

**1st PASS** [0.75]

item counts:

$1 \Rightarrow 2$, $2 \Rightarrow 2$, $3 \Rightarrow 4$, $4 \Rightarrow 2$, $5 \Rightarrow 2$

**2ND PASS** [0.75]

Freq items: 1, 2, 3, 4, 5

BITMAP ① :

| B0 | B1 | B2 |
|---|---|---|
| 1 | 1 | 1 |

BITMAP ② :

| B0 | B1 | B2 |
|---|---|---|
| 1 | 1 | 1 |

| Buckets | Hash ① | Hash ② |
|---|---|---|
| 0 | 4 | 4 |
| 1 | 4 | 3 |
| 2 | 4 | 5 |

Candidate item pairs :

{1,2} {1,3} {1,5} {2,3} {2,4}
{3,4} {3,5} {4,5}

Counts of pairs

{1,2} ⇒ 1    {1,3} ⇒ 2    {1,5} ⇒ 1
{2,3} ⇒ 2    {2,3} ⇒ 1    {3,4} ⇒ 2
{3,5} ⇒ 2    {4,5} ⇒ 1

Frequent Pairs :- {1,3}, {2,3} {3,4}, {3,5}

# MULTI STAGE

## 1st PASS    [0·5 POINTS]

item counts :
$$1 \Rightarrow 2, \quad 2 \Rightarrow 2, \quad 3 \Rightarrow 4, \quad 4 \Rightarrow 2, \quad 5 \Rightarrow 2$$

### HASH ①

| BUCKET | COUNT |
|--------|-------|
| 0 | 4 |
| 1 | 4 |
| 2 | 4 |

## 2ND PASS    [0.5 POINTS]

frequent items : 1, 2, 3, 4, 5

BITMAP ① :

| B0 | B1 | B2 |
|----|----|----|
| 1 | 1 | 1 |

### HASH ②

| BUCKET | COUNT |
|--------|-------|
| 0 | 4 |
| 1 | 3 |
| 2 | 5 |

## 3RD PASS    [0·5 POINTS]

frequent items : 1, 2, 3, 4, 5

BITMAP ②

| B0 | B1 | B2 |
|----|----|----|
| 1 | 1 | 1 |

Count of pairs :

$$\{1,2\} \Rightarrow 1, \quad \{1,3\} = 2, \quad \{1,5\} \Rightarrow 1, \quad \{2,3\} \Rightarrow 2, \quad \{2,4\} \Rightarrow 1,$$
$$\{3,4\} \Rightarrow 2, \quad \{3,5\} \Rightarrow 2, \quad \{4,5\} = 1$$

Frequent Pairs : $\{1,3\}, \{2,3\}, \{3,4\}, \{3,5\}$

2) General idea - [0.5]

for some applications, it is sufficient to discover most frequent itemsets and is not essential every to discover every single one.

Pros: less I/o cost, time,  Cons: False positive, False
[0.25]  Passes.  [0.25] Negative are induced in the results.

3) PHASE 1 MAP i/p :-                     PHASE 1 MAP O/p :-
[0.25] a chunk/subset of all baskets      Set of pairs (f,1) where f is
(sample of i/p file)                       a frequent itemset from sample.

PHASE 1 REDUCE 1/p :-                      PHASE 1 REDUCE O/p :-
[0.25] Set of pairs (f, 1)                 Candidate itemsets

PHASE 2 MAP i/P :-                         PHASE 2 MAP O/p :-
[0.25] Result from phase1 and              Set of pairs (c,v), c is
total output file                          candidate itemsets, v is the
                                           support for that itemset.

PHASE 2 Reduce i/p                         PHASE 2 REDUCE O/P :-
[0.25] . (c,v)                             if v ≥ s, emit (c,v)

4) false positives :- Infrequent in entire data, frequent in
[0.25]                                           Sample.

[0.25] false Negative :- Frequent in entire data, Infrequent in
                                                  Sample.

Increasing support : will increase FN as it will be
harder to be frequent in the sample, decrease FP.
[0.5] Decreasing support : will induce more FP in the data
as it is easier to be frequent in the data/sample

decrease FNs.

**5)**

[0.25] Singleton: $\{a\}$ is in -ve border,

           iff $\{a\}$ is not frequent in the sample

[0.25] Pair: $\{a,b\}$ is in -ve border,

       iff $\{a,b\}$ is not frequent in the sample

        $\{a\}, \{b\}$ are frequent

0.5 marks

① Construct sample data set

② Find candidates frequent itemsels from sample

③ Construct Negative border.

④ Process the whole file,

    if no itemset from -ve border turns out to be frequent in whole data set, correct ✓

    if some ...., Repeat the algo with random sample.