

## Assignment 4

### Question 1:

Only  $N$  of that million items are frequent, so we're only interested in keeping an array of length  $N$  to keep the counts of the frequent items for our second pass. This will take up  $4N$  bytes and will replace the 4 million mentioned above.

The remainder of the memory required will be taken up with our count of the frequent pairs, which we can do using a triangular array or a hash table of triples.

The triangular array will consist of a one dimensional array of length  $(N^2)$  which at 4 bytes per slot, takes up  $4N^2$  bytes. So for the triangular array method, our total number of bytes  $S$ , will be calculated as follows

$$S = 4N^2 + 4N //$$

## Question 2

	B	C	P	M	J
1	x		x		
2		x		x	
3	x	x			x
4			x	x	
5	x	x		x	
6				x	x
7			x		x
8	x	x		x	x

B → Bread

C → Coke

P → Pepsi

M → Milk

J → Juice

$$s(B, C) = 3$$

$$s(B, P) = 1$$

$$s(B, J) = 2$$

$$s(B, M) = 2$$

$$s(C, P) = 0$$

$$s(C, M) = 3$$

$$s(C, J) = 2$$

$$s(P, M) = 1$$

$$s(P, J) = 0$$

$$s(M, J) = 2$$

Support Threshold (s) = 2

So select item which have threshold equal or greater than 2.

$[B, C], [B, M], [B, J]$

$[C, M], [C, J], [M, J]$

### Question 3

Expected memory for pass 2:

Expected no. of frequent buckets = 1,000,000

Each frequent bucket hashes to  $\left[1 + \frac{P}{\text{buckets}}\right]$

pairs which will simplify to  $\frac{P}{\text{bucket}}$  pairs.

Probability of bucket to be frequent is map to frequent bucket  $\Rightarrow \frac{1,000,000}{\text{buckets}}$

Total expected memory consumption for pass 2

$$12 \text{ bytes/pair} = \frac{P \times 12,000,000}{\text{buckets}}$$

In pass 3

Let us use 4 mb space to count item & remainder of S as a hash table to help eliminate non frequent pairs

$$\textcircled{1} \text{ integers (4 bytes)} = \frac{S \sim 4 \text{ MB}}{4} \approx \frac{S}{4} \text{ buckets hash table}$$

This table is compressed to bit map before

pass.

$\textcircled{2}$  But no. of buckets is  $S/4$ . The buckets will have less space in pass 2.

$\therefore S/4 \rightarrow$  upper bound for buckets



$S_1$  buckets are compressed to  $S_1$  bytes using counting

$S = \frac{31}{32}$  bytes for counting problem, on per the assumption, 2<sup>nd</sup> pass

$$\left[ \frac{p = 12,000,000}{\text{buckets}} \right] \text{ bytes for counting so it buckets} = S/4$$

$$\text{we need } p = 12,000,000 / (S/4)$$

$$= \frac{18,000,000 \times 4}{S} \text{ bytes for counting since we}$$

$$\text{have } S = \frac{31}{32} \text{ bytes tree}$$

$$S = \frac{31}{32} = \frac{18,000,000 \times 4}{S}$$

$$S^2 \times 31 = 18,000,000 \times 32 \times 4$$

$$p = \frac{S^2 \times 31}{18,000,000 \times 32}$$

$$\text{Upper bound is } p < \frac{S^2}{19,548,387}$$

### Question 11

Set of item :  $\{A, B, C, D, E, F, G, H\}$

Maximal frequent itemset

$\{A, B\}$   $\{A, C\}$   $\{A, D\}$  ,  $\{B, C\}$  ,  $\{E\}$  ,  $\{F\}$

→ For singleton sets

Singletons sets which are not in frequent sets are in negative border.

$\{G\}$   $\{H\}$  → -ve border

→ For pairs item will be in -ve border when all of its subsets are in frequent subset but not set.

$\{A, E\}$   $\{A, F\}$   $\{B, D\}$   $\{B, E\}$   $\{B, F\}$   $\{C, D\}$   
 $\{C, E\}$   $\{C, F\}$   $\{D, E\}$   $\{D, F\}$   $\{E, F\}$

→ For triplen →  $\{A, B, C\}$

sets in -ve border :  $\{G\}$   $\{H\}$

$\{A, E\}$   $\{A, F\}$   $\{B, D\}$   $\{B, E\}$   $\{B, F\}$   $\{C, D\}$  ,  
 $\{C, E\}$   $\{C, F\}$   $\{D, E\}$   $\{D, F\}$   $\{E, F\}$   $\{A, B, C\}$