

Frequent itemsets

Question 1:

Suppose we have transactions that satisfy the following assumptions:

- s , the support threshold, is 10,000.
- There are one million items, which are represented by the integers $0, 1, \dots, 999999$.
- There are N frequent items, that is, items that occur 10,000 times or more.
- There are one million pairs that occur 10,000 times or more.
- There are $2M$ pairs that occur exactly once. M of these pairs consist of two frequent items, the other M each have at least one nonfrequent item.
- No other pairs occur at all.
- Integers are always represented by 4 bytes.

Suppose we run the a-priori algorithm to find frequent pairs and can choose on the second pass between the triangular-matrix method for counting candidate pairs (a triangular array $\text{count}[i][j]$ that holds an integer count for each pair of items (i, j) where $i < j$) and a hash table of item-item-count triples. Neglect in the first case the space needed to translate between original item numbers and numbers for the frequent items, and in the second case neglect the space needed for the hash table. Assume that item numbers and counts are always 4-byte integers.

As a function of N and M , what is the minimum number of bytes of main memory needed to execute the a-priori algorithm on this data?

----> One data structure is needed to hold the counts of each item. This will be an array of length 1,000,000 which at 4 bytes an integer, is 4 million bytes. Keeping an array of length N will take up $4N$ bytes to keep the counts of the frequent items. A hash table is needed to hold M values. The two items in the pair and the count will be recorded so that $3 \text{ integers} \times 4 \text{ bytes} = 12 \text{ bytes per integer}$, so the size of this will be $12M$.

The minimum number of bytes of main memory needed to execute the a-priori algorithm on this data is $S = 4N + 12M$

Question 2:

Below is a table representing eight transactions and five items: Beer, Coke, Pepsi, Milk, and Juice. The items are represented by their first letters; e.g., "M" = milk. An "x" indicates membership of the item in the transaction.

	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

Compute the support for each of the 10 pairs of items. If the support threshold is 2, find out the pairs that are frequent itemsets.

The pairs of item sets are:

(B, C), (C, M) \rightarrow 3

(B, J), (B, M), (C, J), (M, J) \rightarrow 2

(B, P), (P, J), (P, M) \rightarrow 1

(C, P) \rightarrow 0

The pairs of frequent item sets are (B, J), (B, M), (C, J), (M, J)

Question 3:

Suppose we perform the PCY algorithm to find frequent pairs, with market-basket data meeting the following specifications:

- s , the support threshold, is 10,000.
- There are one million items, which are represented by the integers 0,1,...,999999.
- There are 250,000 frequent items, that is, items that occur 10,000 times or more.
- There are one million pairs that occur 10,000 times or more.
- There are P pairs that occur exactly once and consist of 2 frequent items.
- No other pairs occur at all.
- Integers are always represented by 4 bytes.
- When we hash pairs, they distribute among buckets randomly, but as evenly as possible; i.e., you may assume that each bucket gets exactly its fair share of the P pairs that occur once.

Suppose there are S bytes of main memory. In order to run the PCY algorithm successfully, the

number of buckets must be sufficiently large that most buckets are not frequent. In addition, on the second pass, there must be enough room to count all the candidate pairs. As a function of S , what is the largest value of P for which we can successfully run the PCY algorithm on this data? Find out the value for S and value for P that is approximately (i.e., to within 10%) the largest possible value of P for that S .

Assignment - 4

(30) $S = 1000000$, Items = 1000000

$P(\text{Probability of bucket to be freq}) = \frac{1000000}{\text{buckets}}$

No of pairs that map to be a freq bucket = $P \times \frac{1000000}{\text{buckets}}$

→ In pass 1, we need some space to count items (~4MB) & we have at most $\frac{(S - 4MB)}{H} \approx \frac{S}{H}$ (buckets in this hash table).

→ According to analysis, on the 2nd pass we'll need $(P \times \frac{12000000}{\text{buckets}})$ byte for counting.

So, if buckets = $\frac{S}{H}$.

we'll need $(P \times \frac{12000000}{\frac{S}{H}}) = \frac{480000000 \times P}{S}$ bytes for counting.

Since we have $S \times \frac{31}{32}$ bytes free =

$\Rightarrow S \times \frac{31}{32} = 48,000,000 \times \frac{P}{S}$

$\Rightarrow S^2 = 49,548,387$

$\Rightarrow P = \frac{S^2}{49,548,387}$

bounds:

$P \leq \frac{S^2}{49,548,387}$

Question 4: During a run of Toivonen's Algorithm with set of items {A,B,C,D,E,F,G,H} a sample is found to have the following maximal frequent itemsets: {A,B}, {A,C}, {A,D}, {B,C}, {E}, {F}. Compute the negative border.

---->The negative border consists of fourteen sets: {G}, {H}, {A, E}, {A, F}, {B, D}, {B, E}, {B, F}, {G, D}, {C, E}, {C, F}, {D, E}, {D, F}, {E, F}, {A, B, C}