

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,...,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?

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

Explanation- One data structure is needed to hold the count of each item. This will be an array of 4 million bytes. If the length of the array is N , then it takes upto $4N$ bytes to store the count of frequent items. A hash table is needed to store M values. The two items and count will be stored, so the size will be $12M$.

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 pair of itemsets-

- (B,C),(C,M)-->3
- (B,J),(B,M),(C,J),(M,J)-->2
- (B,P),(P,J),(P,M)-->1
- (C,P)-->0

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 .

$S=10000$, Items=1000000, $P=1000000/\text{buckets}$

No of frequent pairs that map to a bucket= $p*(1000000/\text{buckets})$

During pass1, we have utmost $(S-4MB)/4$ approx $S/4$ buckets

For pass 2, we need $P*12000000/\text{buckets}$

To have enough space for all the counts, we need $S \geq 48000000P/S$ or

$P \leq S^2/48000000$

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-

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