

Week 2: Frequent Itemsets (Advanced)

[Help Center](#)

Warning: The hard deadline has passed. You can attempt it, but **you will not get credit for it**. You are welcome to try it as a learning exercise.

☐ In accordance with the Coursera Honor Code, I (Manuel Bordés Rguez.) certify that the answers here are my own work.

Question 1

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 $0, 1, \dots, 999999$.
- 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 large. 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? Demonstrate that you have the correct formula by indicating which of the following is a value for S and a value for P that is approximately (i.e., to within 10%) the largest possible value of P for that S .

- ☐ $S = 1,000,000,000$; $P = 35,000,000,000$
- ☐ $S = 500,000,000$; $P = 5,000,000,000$

- ☐ S = 100,000,000; P = 120,000,000
- ☐ S = 200,000,000; P = 400,000,000

Question 2

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. Then, identify in the list below the set that is NOT in the negative border.

- ☐ {H}
- ☐ {A,B,D}
- ☐ {B,E}
- ☐ {G}

☐ In accordance with the Coursera Honor Code, I (Manuel Bordés Rguez.) certify that the answers here are my own work.

Submit Answers

Save Answers

You cannot submit your work until you agree to the Honor Code. Thanks!

