



COMP-4250 Big Data Analytics and Database Design

Project I (15%) Mining Frequent Itemsets

Deadline: End of Friday February 14th 2020

Important Note: This project can be done in a **group of two** or **individually**. If you want to do the project in a group of two, you have to send the name of your teammate to the instructor via email no later than end of Friday January 31st, and otherwise it will be assumed that you perform the project individually.

Description

The main objective of this project is to **find frequent itemsets by implementing two efficient algorithms**: **A-Priori** and **PCY**. The goal is to **find frequent pairs** of elements. You do not need to find triples and larger itemsets.

Resources

Lectures 2 and 3 on Blackboard, and Chapter 6 of the textbook.

Programming Language

You can choose your favorite programming language, preferably one of the following ones: C, C++, Java, C#, or Python.

Dataset

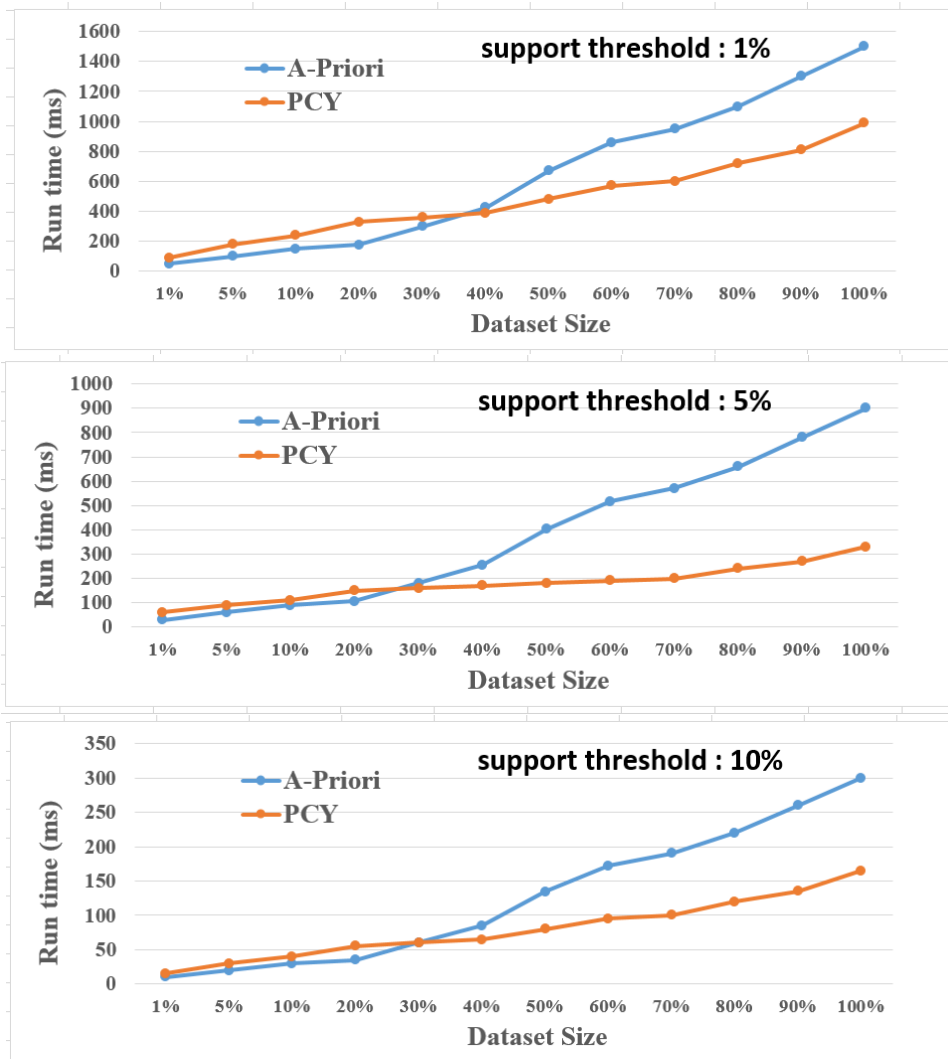
The retail dataset contains anonymized retail market basket data (88K baskets) from an anonymous retail store. The preprocessing step to map text labels into integers has already been done. Use Sublime Text, TextPad or Notepad++ or other software to open the file. **Do not use Notepad.**

Dataset link: It is available on course page on Blackboard.

Experiments

Perform the **scalability study** for finding frequent pairs of elements by dividing the dataset into different chunks and measure the time performance. Provide the line chart. Provide results for the following support thresholds: 1%, 5%, 10%. For example, if your chunk is 10% of the dataset, you have around 8,800 baskets. Therefore, if your support threshold is 5%, you should count the pairs that appear in at least 440 baskets. See three samples below for three different support thresholds.

Note: the sample charts contain hypothetical numbers!



Optional (Bonus Points)

- Implement Multistage (3 Passes) version of PCY, using one extra hashtable (0.25% extra). (add the results to the line chart)
- Implement Multihash version of PCY, using one extra hashtable (0.25% extra). (add the results to the line chart)

Submission

You have to submit your **code**, along with a report of your **experiments** on Blackboard before the deadline. Indicate the specification of the machine that you run the experiments on, including the operating system, CPU, and RAM.