# TD3: Declarative Itemset Mining

Pablo Mollá Chárlez

February 4, 2025

# Contents

# 1 Exercise 1

In this exercise, you will work with the following tools:

- **Choco-Mining:** A Java library designed for solving itemset mining problems, built on the Choco-solver framework.

- **The SPMF library:** An open-source Java-based software and data mining library specializing in pattern mining (SPMF).

## 1.1 Question 1
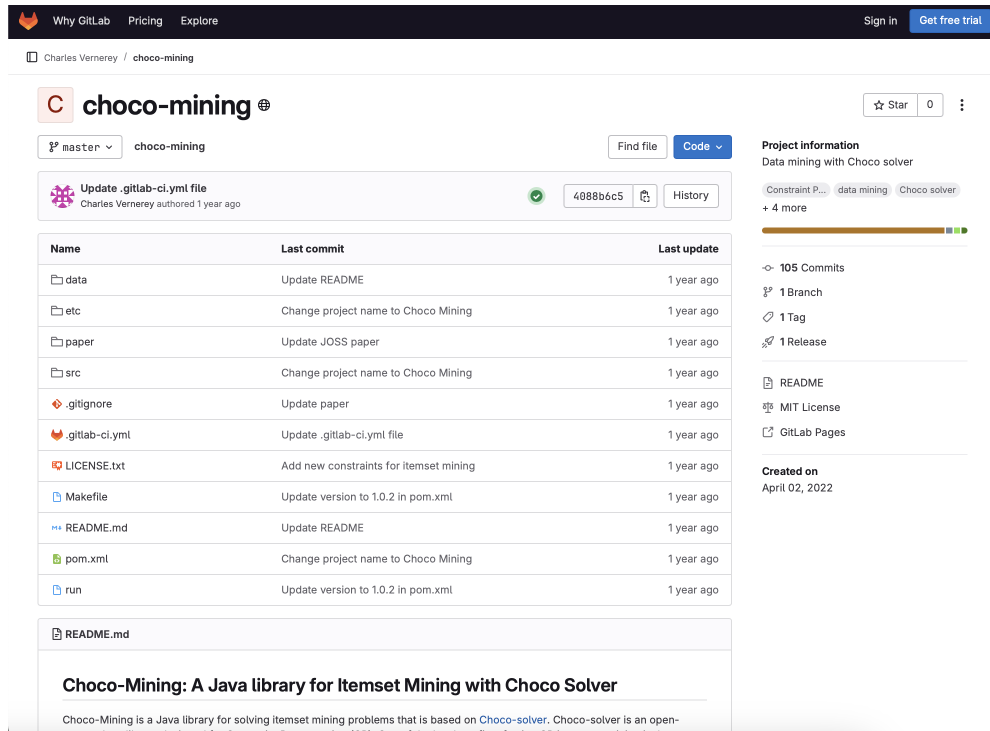
Clone the GitHub repository of Choco-Mining (link)



Figure 1: Choco-Mining Repository

## 1.2 Question 2

Open the file ExampleClosedItemsetMining.java and perform the following tasks:

1. Review the code in detail.

2. Run the main method.



Figure 2: Main Method Result

3. Run it on other datasets such as mushroom or chess.



Figure 3: Mushroom Modification Code



Figure 4: Chess Modification Code

4. Display the number of resulting patterns.

5. Display the execution time.



Figure 5: Patterns and Time Execution Modification

Figure 6: Patterns and Time Execution for Mushroom Dataset

## 1.3 Question 3

Add the frequency constraint: $freq(P) \geq \alpha$

**In next section, the image answers the question.**

## 1.4 Question 4

Add a constraint on the size of the returned patterns: $size(P) \geq lb$.

```
int alpha = 10;
IntVar freq = model.intVar( name: "freq", alpha, database.getNbTransactions());
// Integer variable that represents the length of x with the bounds [1, nbItems]
int lb = 10;
IntVar length = model.intVar( name: "length", lb, database.getNbItems());
```

Figure 7: Constraint on frequency and size of returned patterns

## 1.5 Question 5

Now, replicate the tasks using SPMF. Run the **.jar** file available in your local repository. The goal is to run LCM for closed itemset enumeration, relaunch with different thresholds for frequency, and also for pattern size.

```
(base) chenchenjunjie@client-172-18-83-132 DIM-20250203 % java -jar spmf.jar run LCM choco-mining/data/mushroom.dat output.txt 40.0%
>/Users/chenchenjunjie/m2/dm_cp/DIM-20250203/spmf.jar
========== LCM - STATS ============
 Freq. closed itemsets count: 107
 Total time ~: 81 ms
 Max memory:51.803619384765625
==================================
```

Figure 8: SPMF: Mushroom Dataset with 40% Frequency

```
(base) chenchenjunjie@client-172-18-83-132 DIM-20250203 % java -jar spmf.jar run LCM choco-mining/data/mushroom.dat output.txt 50.0%
>/Users/chenchenjunjie/m2/dm_cp/DIM-20250203/spmf.jar
========== LCM - STATS ============
 Freq. closed itemsets count: 44
 Total time ~: 49 ms
 Max memory:31.412551879882812
--------------------------------------
```

Figure 9: SPMF: Mushroom Dataset with 50% Frequency

```
[(base) chenchenjunjie@client-172-18-83-132 DIM-20250203 % java -jar spmf.jar run LCM choco-mining/data/mushroom.dat output.txt 70.0%
>/Users/chenchenjunjie/m2/dm_cp/DIM-20250203/spmf.jar
========== LCM - STATS ============
 Freq. closed itemsets count: 12
 Total time ~: 28 ms
 Max memory:17.77068328857422
==================================
```

Figure 10: SPMF: Mushroom Dataset with 70% Frequency

Now, we apply modifications at both levels, **frequency** and **pattern size**.



Figure 11: SPMF: Mushroom Dataset with 70% Frequency and different pattern sizes



Figure 12: SPMF: Mushroom Dataset with different frequencies and pattern size (= 5)

### 1.6  Question 6

Add a constraint, called CategoryConstraint, to the file **ExampleClosedItem-setMining.java** to model the following problem : Consider a dataset with n items, organized into categories of size **catSize** (e.g., household products, appliances, etc.). The dataset is divided into $nbCat = \frac{n}{catSize}$ categories, with items that do not belong to any category (but do not exceed the size of **catSize**). Figure **??** shows an example with 8 items, 2 categories of size 3, and 2 items that do not belong to any category. The task is to create a constraint model that enumerates all closed itemsets composed of items belonging to at least m categories:

$$\text{CategoryConstraint(P)} = \sum_{i=1}^{nbCat} \prod_{j=1}^{catSize} P_i \geq m$$

For example, in the dataset shown in Figure **??**, with $m = 2$, the following pattern is produced: **BEF**.

### 1.7  Question 7

How can this **CategoryConstraint** be taken into account in SPMF ?