



# CSCI 5523 Introduction to Data Mining (sec 001, 771, 881, 883, 885) Fall 2015

My Home ► CSCI5523\_001F15 ► General ► Project #1 (due 10/19/2015 @ 11:55pm)

## Project #1 (due 10/19/2015 @ 11:55pm)

**[Updated on 10/15/2015]** In order to simplify the assignment, do not generate association rules for support values of 20 and 15 but you need to generate them for the other support values. However, you do need to generate the frequent itemsets for all support values. As a result of this change, you can assume that you can store the frequent itemsets in memory so that you can use them during association rule generation.

### Objective

The purpose of the project is to write a program to generate all association rules whose support is greater than a user-supplied minimum support and whose confidence is greater than a user supplied minimum confidence. You need to implement the recursive database-projection based algorithm that I described in class and is included in the lecture slides. This class of algorithms are also describe in Charu's textbook in section 4.4.3.2 (though the description there is more general than the method that I presented in class).

Your program should take as command line option five parameters: (i) the minimum support, (ii) the minimum confidence, (iii) the name of the input file, (iiii) the name of the output file, and (iv) options. The specific parameter sequence along with the name of the executable are as follows:

```
hcrminer minsup minconf inputfile outputfile options
```

NOTE: For those using Java, the command line should look like the following:

```
java hcrminer minsup minconf inputfile outputfile options
```

The options parameter will be a numerical value whose meaning is as follows:

options = 1

The numbering of the items coming from the input file is used as the lexicographical ordering of the items.

options = 2

The lexicographical ordering of the items is determined by sorting the items in increasing frequency order in each projected database.

options = 3

The lexicographical ordering of the items is determined by sorting the items in decreasing frequency order in each projected database.

Your program must be written in either C, C++, or Java.

## Input file format

The input file consists of a set of lines, each line containing two numbers. The first number is the transaction ID and the second number is the item ID. The lines in the file are ordered in increasing transaction ID order. Note that the set of items that make up the transaction will be derived by combining the item IDs of all the lines that correspond to the same transaction ID.

Two input files are provided. The first is a small one that you can use during initial code development. The second is larger and will be the one on which you need to report performance results.

## Output file format

The output file will contain as many lines as the number of high-confidence frequent rules that you found. The format of each line will be as follows:

```
LHS | RHS | SUPPORT | CONFIDENCE
```

Both LHS and RHS will contain the items that make up the left- and right-hand side of the rule in a space delimited fashion.

For minsup of 15 and 20, the format of each line should be:

```
LHS | { } | SUPPORT | -1
```

where LHS is the frequent itemset that has minsup of 15 or 20.

## Report

Along with your code, you need to submit a report that contains the following:

A description of the data structures that you use to store the projected database and how you implemented the projection.

Run your code for each combination of the following sets of values for the minimum support, the minimum confidence, and option:

```
minsup: {15, 20, 30, 50, 100, 500, 1000} [Note that this is the support count (actual frequency)]
minconf: {0.8, 0.9, 0.95}
options: {1, 2, 3}
```

For each combination of minconf and options value, generate two bar-charts. The first will show the amount of time required to find the frequent itemsets and to find the rules (a set of bars for each of these two quantities) for the different values of minimum support. The second will show the number of frequent itemsets that were found and the number of high confidence rules that were generated (a set of bars for each of these two counts) for the different values of minimum support. Make sure that each bar is also annotated with the quantity that is plotting (at the top of the bar).

Here is a table of the approximate amount of time taken for a couple of minsup values:

minsup	time (secs)
500	0.4
30	11.04

A discussion as to which value of the options parameter resulted in faster execution time and why.

## Submission

DO FOLLOW ALL OF THESE INSTRUCTIONS

1. You need to submit a directory that contains your code along with the report.
2. The report should be named "report.pdf" and should be in PDF format.
3. Your code should be in a single file named "hcrminer.<extension>".
4. The directory should contain a Makefile such that by simply doing a "make" it will build the "hcrminer" executable.
5. The directory should be uploaded in as a "<umn-username>.tar.gz" archive.


## Grading criteria

1. Report: 10%
2. Correctness: 60%
3. Efficiency: 30%

 large.bz2

 small.bz2

## Submission status

Submission status	Submitted for grading
Grading status	Not graded
Due date	Monday, October 19, 2015, 11:55 PM
Time remaining	Assignment was submitted 3 hours 56 mins early
Last modified	Monday, October 19, 2015, 7:59 PM
File submissions	 umn-sadal005.tar.gz
Submission comments	► Comments (0)

## ADMINISTRATION



Course administration