

# Projekt MED-P3, algorytm GRM. Raport.

Przedmiot: Metody eksploracji danych w odkrywaniu wiedzy.

Michał Aniserowicz, Jakub Turek

## 1 Opis zadania

Celem projektu jest zaimplementowanie algorytmu wyznaczania reguł decyzyjnych o minimalnych poprzednikach, które są częstymi generatorami. Algorytm ten jest modyfikacją algorytmu odkrywania częstych generatorów (GRM), opisanego w [1].

## 2 Założenia

poczynione zalozenia - kazda transakcja bedzie miala decyzje - aplikacja konsolowa pobierajaca dane z pliku i zwracajaca wynik w dwóch formatach - aplikacja mierzy czas poszczególnych kroków - C#, .NET 3,5

## 3 Dane wejściowe i wyjściowe

opis danych wejściowych i wyjściowych - opcje (zostana opisane pozniej), minsup bezwzględne! - dane oddzielone przecinkami (decyzja razem z atrybutami, na dowolnym miejscu) - nagłówki w pierwszym wierszu - brak danych - spacja (białe znaki) - dwa formaty wyników - oprócz tego wynik na konsoli

## 4 Implementacja

wszystkie istotne kwestie związane z projektowaniem (np. diagramy klas) i implementacja projektowanie:  
- podział na moduły (console, dataset processing, GRM) - testy - diagram klas Logic implementacja:  
- jakiś algorytm, może z diffsetami - różne sortowania - tidset/diffset - bruteforce/inv list - tracking (poziomy)

### 4.1 Optymalizacje

- wszystkie wartości otrzymują identyfikatory liczbowe - skonfliktowane generatory - transaction ids - posortowane (szybkie intersect, except)  
    różnice z GRM: - dany node jest decyzyjny - nie rozwijamy go (bo generatory dzieci nie będą minimalne) - generatory decyzji trzymane w słowniku (klucz - decyzja), posortowane wg hashu - w ogóle nie ma granicy - dla diffsetów transaction ids trzymane w słowniku (klucz - decyzja)

## 5 Podręcznik użytkownika

podręcznik potencjalnego użytkownika wytworzonego oprogramowania (zamierzam korzystać z niego podczas sprawdzania Państwa rozwiązań) - wszystkie opcje programu - przykładowa komenda i wynik na konsoli

## 6 Analiza poprawności

wszystkie wyniki wytwarzane przez program otrzymane dla małego, przykładowego zbioru danych (w celu weryfikacji poprawności działania programu) - przykład z konsultacji

## 7 Analiza wydajności

wyniki jakościowe i ilościowe na (np. czas działania; liczba wzorców) uzyskane dla większych (wielkich) zbiorów danych (np. z <http://archive.ics.uci.edu/ml/> or <http://fimi.cs.helsinki.fi/data/> lub uzgodnionych już wcześniej ze mną podczas konsultacji projektowych) - wykresy, wykresy - ze dla dużej liczby atrybutów mało wydajny

## 8 Wnioski

wnioski z realizacji projektu - ze trzeba by poprawić wykrywanie supergeneratorów - ze ogólnie działa spoczko (nurse)

## Literatura

- [1] *Odkrywanie reprezentacji generatorowej wzorców częstych z wykorzystaniem struktur listowych*, Kryszkiewicz M., Pielasa P., Instytut Informatyki, Politechnika Warszawska.