

Peti domaći zadatak iz Tehnika i metoda analize podataka

Danilo Veljović, broj indeksa 1120

Opis skupa podataka

Skup podataka korišćen za testiranje implementiranog algoritma je skup podataka o sertifikovanim biznisima sa Small Business Services sajta. Dataset sadrži informacije o tome ko su vlasnici biznisa, tj kog su pola, kog su manjinskog ili većinskog porekla, na kojoj lokaciji se nalaze i koliko ima takvih biznisa na toj lokaciji.

Dataset sadrži 1147 instanci.

Način implementacije algoritma

Ključni deo algoritma je dat u funkciji runApriori() od linije 81 do linije 89. Najpre se generiše kandidat skup, nadovezivanjem podskupa iz prethodnog koraka. Zatim se izvodi testiranje i pruning tog skupa, ako se tu nalazi neki ne-česti (infrequent) podskup. Na kraju se računa najčešći skup pribavljanjem onih elemenata koji zadovoljavaju minimalnu potporu (minimal support). Slika 1 ima isečak iz koda koji radi upravo ovo (u kodu linije 82 - 89)

```
largeSet[k - 1] = currentLSet
currentLSet = joinSet(currentLSet, k)
currentCSet = returnItemsWithMinSupport(
    currentLSet, transactionList, minSupport, freqSet
)
```

Slika 1 – Srž Apriori algoritma

Generisanje itemseta (skupova) se radi funkcijom joinSet() kojom se kreira unija svih elementa koji su određene dužine.

Funkcijom returnItemsWithMinSupport() se određuje i vraća podskup elementa skupa čija je potpora veća od unapred zadate minimalne potpore. Potpora se računa formulom na slici 2.

$$support(I) = \frac{\text{Number of transactions containing } I}{\text{Total number of transactions}}$$

Na kraju algoritma se finalni skup elemenata još jednom filtrira tako da se sad računa poverenje za svaki element podskupa i ako je vrednost nekog elementa veća od minimalne zadate vrednosti, taj element će biti vraćen. Na slici 3 je data izlazna vrednost algoritma.

```

C:\Users\danil\Desktop\Apriori-python3>python apriori.py -f INTEGRATED-DATASET.csv
item: ('Brooklyn',) , 0.152
item: ('HISPANIC',) , 0.164
item: ('HISPANIC', 'MBE') , 0.164
item: ('MBE', 'WBE') , 0.169
item: ('MBE', 'New York') , 0.170
item: ('WBE', 'New York') , 0.175
item: ('MBE', 'ASIAN') , 0.200
item: ('ASIAN',) , 0.202
item: ('New York',) , 0.295
item: ('NON-MINORITY',) , 0.300
item: ('NON-MINORITY', 'WBE') , 0.300
item: ('BLACK',) , 0.301
item: ('MBE', 'BLACK') , 0.301
item: ('WBE',) , 0.477
item: ('MBE',) , 0.671

----- RULES:
Rule: ('WBE',) ==> ('NON-MINORITY',) , 0.628
Rule: ('ASIAN',) ==> ('MBE',) , 0.990
Rule: ('HISPANIC',) ==> ('MBE',) , 1.000
Rule: ('BLACK',) ==> ('MBE',) , 1.000
Rule: ('NON-MINORITY',) ==> ('WBE',) , 1.000

```

Slika 3 – Rezultat izvršenja Apriori algoritma

Prvi elementi izlaza su vezani za sve elemente izlaznog skupa i odgovarajuće vrednosti njihovig potpora.

Drugi deo rezultata (od RULES podnaslova) vezan je direktno za pravila koja su zaključena i za njihov nivo poverenja.

U programu se može zadati i minimalni nivo poverenja, kao minimalni nivo potpora. Defaultne vrednosti su *minSupport* = 0.15 i *minConfidence* = 0.6.