

Eksploracja dużych wolumenów danych

Projekt – Sprawozdanie I

Sformułowanie problemu i przetwarzanie danych

ImięTM Nazwisko ImięTM Nazwisko

27 listopada 2012

1 Sformułowanie problemu

Krótką informacją o wybranym konkursie: nazwa, strona internetowa, nagroda :), założyciel.

1.1 Krótka opis problemu

Krótki i konkretny opis zadania konkursowego.
(maks. 1/2 strony)

1.2 Opis danych

Opis zbioru uczących, walidujących (jeżeli jest) i testowego wraz z podstawowymi statystykami: rozmiar, liczba atrybutów, liczba przykładów.

Warto przedstawić postać typowych przykładów uczących oraz testowych.

Należy również opisać sposób ostatecznej weryfikacji algorytmów: trwający konkurs (wynik końcowy w konkursie; ocena na podstawie tablicy wyników, czyli na zbiorze walidującym), możliwość weryfikacji po zakończeniu konkursu na jego stronie, wewnętrzny zbiór testowy (rozmiar, jak zostanie wyodrębniony, kiedy zostanie wykorzystany, itp.).
(maks. 3/4 strony)

1.3 Opis metody oceny rozwiązania,

Dokładny i zwięzły opis metody oceny rozwiązania (jeżeli algorytm będzie testowany wewnętrznie, to należy napisać odpowiednie oprogramowanie oceniające rozwiązanie).

(maks. 1/2 strony)

2 Reprezentacja danych w pamięci zewnętrznej

Krótki opis (bardzo zwięzły i konkretny) z uwzględnieniem:

- opisu reprezentacji z jej zaletami i wadami, również z perspektywy eksperymentalnej, które należy przeprowadzić,
- dyskusji na temat konfiguracji i administracji danego rozwiązania,
- dyskusji na temat czasu utworzenia i objętości reprezentacji danych,
- dyskusji na temat napotkanych trudności.

Można opisać więcej niż jedno rozwiązanie, w celach porównawczych.
(maks. 2 strony)

3 Reprezentacja danych w pamięci operacyjnej

Krótki opis (bardzo zwięzły i konkretny) z uwzględnieniem:

- opisu reprezentacji z jej zaletami i wadami, również z perspektywy eksperymentalnej, które należy przeprowadzić,
- dyskusji na temat czasu wczytywania do pamięci, dostępu do danych i objętości danych,
- dyskusji na temat napotkanych trudności.

Można opisać więcej niż jedno rozwiązanie, w celach porównawczych.
(maks. 2 strony)

4 Eksperyment

4.0.1 Wyszukiwanie najbliższych sąsiadów

Wyszukiwanie najbliższych sąsiadów zostało przetestowane dla 9, 86, 862 i 8623 lekcji oraz x, xx, xxx i xxxx autorów czyli dla odpowiednio 0,1%, 1%, 10% i 100%. Wyniki czasowe wraz z odchyleniem standardowym zebrano w tabelach.

W tym rozdziale należy umieścić wyniki eksperymentów. Należy zamieścić krótkie komentarze oraz tabele z wynikami (Tabela 2 jest przykładem jak stworzyć tabelę w \LaTeX u).

Należy wykonać następujące zapytania:

1. Grupowanie,
2. Wyszukiwanie najbliższych $k = 50$ sąsiadów.

Ilość lekcji	Procent lekcji	Średni czas odpowiedzi	Odchylenie standardowe
9	0,1%	6754,6	734,93
86	1%	56039,4	4509,16
862	10%	6754,6	734,93
8623	100%	6754,6	734,93

Tablica 1: Wyniki eksperymentu dla wideo lekcji

To są kolumny	zawierające np.	czas odpowiedzi na zapytanie itp.
Opis i liczby:	99.99%	99.99%
Opis i liczby:	99.99%	99.99%
Opis i liczby:	99.99%	99.99%

Tablica 2: Wyniki eksperymentalne

Dla otrzymanych partycji danych (przez grupowanie i wyszukanie najbliższych sąsiadów) należy zastosować dowolną funkcję agregującą na dowolnych atrybutach. Oczywiście, najlepiej jest wybrać atrybut decyzyjny (wyjściowy).

Dla wszystkich eksperymentów należy podać czas wykonania.

Najbliższych sąsiadów należy wyszukać dla maksymalnie pierwszych 10 tys. przykładów, ale szukanie odbywa się w całym zbiorze danych.

Cały eksperyment można powtórzyć parokrotnie i uśrednić wyniki (można podać błęd standardowy). Przy wyznaczaniu czasów nie należy brać pod uwagę czasu wyświetlania wyników.

W zależności od rodzaju problemu, proszę, stosować się do poniżej opisanych wytycznych. Jeżeli dany problem nie pasuje do żadnego wyżej opisanego typu, to proszę o kontakt.

(maks. 4 strony (razem z tabelami)).

4.1 Tabelaryczne dane nominalne

W przypadku danych nominalnych należy przeprowadzić grupowanie:

- po wszystkich atrybutach (w przypadku dużej liczby atrybutów, grupowanie powinno dotyczyć tylko pierwszych 20 atrybutów),
- dla każdego z atrybutu z osobna,
- oraz dla każdej pary atrybutów (maksymalnie dla pierwszych 100 atrybutów).

Jako miarę podobieństwa przy wyszukiwaniu najbliższych sąsiadów należy zastosować współczynnik Jaccarda. Mierzy on podobieństwo

między dwoma zbiorami i jest zdefiniowany jako iloraz mocy cząstki wspólnej zbiorów i mocy sumy tych zbiorów:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

gdzie A i B są... zbiorami, które odpowiadają... obiektom opisanym przez atrybuty nominalne (proszę zwrócić uwagę, że atrybut binarny cząstko wskazuje, czy dana wartość występuje, czy też nie).

Dla poniższego przykładu zbiór $A = \{\text{czerwony, duży, kombi}\}$ odpowiada pierwszemu obiektowi, a zbiór $B = \{\text{niebieski, duży, coupe, kabriolet}\}$ drugiemu. Współczynnik Jaccarda wynosi:

$$J(A, B) = \frac{1}{6}$$

kolor	rozmiar	kombi	sedan	coupe	kabriolet	limuzyna	minivan
czerwony	duży	1	0	0	0	0	0
niebieski	duży	0	0	0	1	1	0

Można też dokonać wcześniej binaryzacji wszystkich wielowartościowych atrybutów nominalnych i dalej operować na atrybutach binarnych, które wskazują... czy dany element występuje w zbiorze.

4.2 Tabelaryczne dane numeryczne

W celu pogrupowania danych numerycznych dla każdego atrybutu należy przeprowadzić dyskretyzację wartości numerycznych do 5 wartości. Każda wartość odpowiada przedziałowi wartości dla atrybutu numerycznego. Należy przeprowadzić dyskretyzację według równej cząstki, czyli otrzymane przedziały są... mniej więcej równoliczne.

Dla tak przekształconych atrybutów należy wykonać normalne grupowanie, zgodnie z opisem w punkcie 4.1.

Przy poszukiwaniu najbliższych sąsiadów należy się posłużyć miarą... euklidesową... Wcześniej należy jednak ustandaryzować wartości atrybutów, tzn. od wartości każdego atrybutu odejmujemy średnią... i dzielimy przez odchylenie standardowe (czyli przeprowadzamy zmienne o różnicach jednostkach do zmiennych niemianowanych).

4.3 Tabelaryczne dane nominalne i numeryczne

W przypadku wymieszanych danych nominalnych i numerycznych dzielimy dane na dwie części i wykonujemy eksperymenty osobno dla każdej części zgodnie z opisem w dwóch powyższych punktach.

4.4 Dane tekstowe

Dane tekstowe należy zamienić na dane nominalne, gdzie każdy atrybut odpowiada jednemu słowu. Eksperyment należy wykonać dla 100 najczęściej występujących słów. Można (nie jest konieczne) wykonać odpowiednie przekształcenia danych tekstowych, takie jak stematyzacja lub lematyzacja.

Dla tak przekształconych danych należy postąpić zgodnie z opisem w punkcie 4.1.

4.5 Dane macierzowe

W niektórych konkursach dane mają postać zależności pomiędzy dwoma typami obiektów, np. użytkownikami i produktami. W takim przypadku należy wykonać grupowanie dla każdego typu obiektów osobno. Jeżeli jedyną informacją o obiekcie jest jego identyfikator, to grupowanie dotyczy będzie tylko jednego atrybutu. W innym przypadku należy stosować się do opisów w powyższych punktach.

Wyszukiwanie najbliższych sąsiadów w tym przypadku powinno polegać na wyszukiwaniu podobnych użytkowników ze względu na wybrane produkty, oraz wyszukiwaniu podobnych produktów ze względu na “wybieranych” użytkowników. W obydwóch przypadkach należy wykorzystać współczynnik Jaccarda w celu obliczenia podobieństwa.

5 Podsumowanie

Na końcu jest zawsze miejsce na krótkie podsumowanie (maks. 1/2 strony). **Całość raportu nie może przekraczać 8 stron.**