

Politechnika Poznańska
Wydział Informatyki i Zarządzania
Instytut Informatyki

Praca dyplomowa magisterska

**RÓWNOLEGŁE ODKRYWANIE REGUŁ ASOCJACYJNYH
ZAIMPLEMENTOWANE NA PROCESORY GRAFICZNE**

inż. Tomasz Kujawa

Promotor
dr inż. Witold Andrzejwski

Poznań, 2011

Tutaj przychodzi karta pracy dyplomowej;
oryginał wstawiamy do wersji dla archiwum PP, w pozostałych kopiach wstawiamy ksero.

Spis treści

1	Wstęp	1
1.1	Cel i zakres pracy	2
2	Podstawy teoretyczne	3
2.1	Definicje	3
2.1.1	Model teoretyczny	3
2.2	Aktualna wiedza	4
2.2.1	Algorytm Apriori	4
	Generowanie zbiorów częstych	5
	Generowanie reguł asocjacyjnych	6
2.2.2	Algorytm FP-growth	7
	Literatura	8

Rozdział 1

Wstęp

Proces informatyzacji przedsiębiorstw, rozpoczęty kilka dekad temu, wprowadził światową gospodarkę na nowe, dotąd nieznane tory rozwoju. Skrócenie procesu produkcyjnego, wprowadzenie kontroli komputerowych, czy też skomputeryzowanych maszyn skróciło i ułatwiło produkcję, a także zarządzanie procesami w firmach i przedsiębiorstwach. Przed ludźmi stanęły możliwości, ale także wyzwania, z którymi nigdy wcześniej nikt nie musiał sobie radzić. Zmiany, jakie nastąpiły przez ostatnie trzy dekady są nieodwracalne i zmuszają informatyków do tworzenia nowych aplikacji, które będą w stanie sprostać wymaganiom im stawianym.

Informatyzacja firm, instytucji oraz innych jednostek organizacyjnych powinna realizować dwa podstawowe cele. Z jednej strony powinna ona usprawniać pracę pojedynczego pracownika poprzez automatyzację realizowanych przez niego rutynowych zadań. Dzięki wykorzystaniu możliwości komputerów działania te powinny być wykonywane szybciej i w sposób bardziej niezawodny. Z drugiej strony celem informatyzacji jest wpływanie na działanie całych firm w wyniku wspomagania decyzji kadry zarządzającej przedsiębiorstwami. Szybka analiza bazująca na pełnej i aktualnej informacji o stanie firmy może ułatwić kadrze zarządzającej podejmowanie trafnych i szybkich decyzje o strategicznym znaczeniu dla rozwoju danego przedsiębiorstwa.

Wprowadzenie komputerów do właściwie każdej przestrzeni ludzkiego życia wpłynęło na wyprodukowanie olbrzymich ilości danych. Reprezentowane są one w sposób umożliwiający ich składowanie i przetwarzanie komputerowe przez aplikacje analityczne. W chwili obecnej ludzkość jest świadkiem eksplozji ilości danych produkowanych przez różnego rodzaju systemy komputerowe. Analiza tych danych przynieść może wymierne korzyści nie tylko w kwestiach finansowych, ale również poznawczych. Dzięki analizie danych zebranych w przeszłości możliwe jest lepsze dopasowanie planów w przyszłości - na tej podstawie planowane mogą być np. akcje marketingowe, czy też promocje w supermarketach spożywczych. Wykorzystanie wiedzy uzyskanej w ten sposób jest niezwykle szerokie i może być użyte w każdym rejonie działalności firmy.

Odkrycie zależności pomiędzy zgromadzonymi danymi bez zastosowania narzędzi informatycznych jest procesem bardzo skomplikowanym i wymagającym do realizacji dużo czasu. Przy obecnej złożoności większości systemów oraz rozmiarom danych produkowanych przez te systemy, koszt czasowy jest na tyle duży, że ręczna analiza tych danych stała się niemożliwa. Dlatego też tworzone są narzędzia umożliwiające odkrywanie prawidłowości w dużych zbiorach danych, by człowiek na tej podstawie mógł podejmować decyzje i wyciągać wnioski.

Dział informatyki, który zajmuje się odkrywaniem ukrytych dla człowieka prawidłowości i reguł w danych nazywa się eksploracją danych (ang. *Data Mining*), który jest jednym z etapów procesu *odkrywania wiedzy z baz danych (KDD, ang. Knowledge Discovery in Databases)*. Proces odkrywania wiedzy w bazach danych obejmują zwykle działania bardziej złożone niż tylko eksploracja

danych. Eksploracja danych to proces odkrywania wiedzy w postaci nowych, użytecznych, poprawnych i zrozumiałych wzorców w bardzo dużych wolumenach danych [FPSSU96]. Możliwości stosowania technik eksploracji danych w praktyce, wymagają efektywnych metod przeszukiwania ogromnych plików lub baz danych. Warto przy tym wspomnieć, że tego typu technologie nie są w chwili obecnej dobrze zintegrowane z systemami zarządzania bazami danych.

Eksploracja danych (w literaturze spotkać można również określenie drążenie danych, ekstrakcja danych, pozyskiwanie wiedzy, czy też wydobywanie danych [EN05]) odbywa się najczęściej w środowisku baz lub hurtowni danych, które stanowią doskonałe źródła danych do analizy - głównie ze względu na łatwość dostępu oraz usystematyzowaną strukturę przechowywanych informacji. Ponieważ liczba odkrytych wzorców w wielu przypadkach może być bardzo duża, odkryte wzorce bardzo często zapisuje się w osobnych relacjach bazy lub hurtowni danych. Pozwala to na ich dalsze przetwarzanie w trybie off-line przez użytkowników końcowych. Pojęcie eksploracji zyskuje coraz większą popularność (również w wymiarze marketingowym) i jest wykorzystywane w wielu dziedzinach ludzkiego życia.

Jednym z najczęściej wykorzystywanych modeli wiedzy w eksploracji danych są reguły asocjacyjne. Reguła asocjacyjna ma postać $X \Rightarrow Y$, gdzie X oraz Y są wzajemnie rozłącznymi zbiorami elementów. Przykładem reguły, która mogła zostać odkryta w badzie danych sklepu komputerowego, może być reguła postaci *komputer* \wedge *myszka* \Rightarrow *monitor*. Reguła ta prezentuje fakt, że klienci kupujący komputer oraz myszkę z dużym prawdopodobieństwem kupią również monitor. W [Agr94] po raz pierwszy sformułowany został problem odkrywania reguł asocjacyjnych. Podstawą wielu algorytmów odkrywania reguł asocjacyjnych jest algorytm Apriori, zaprezentowany w [AIS93].

W ostatnich latach pojawiły się nowe możliwości wykorzystania współczesnych komputerów. W roku 2007 firma NVIDIA udostępniła programistom uniwersalną architekturę obliczeniową CUDA (ang. *Compute Unified Device Architecture*), który umożliwia wykorzystanie mocy obliczeniowej procesorów graficznych (*GPU*, ang. *Graphics Processing Unit*), bądź innych procesorów wielordzeniowych, do rozwiązywania ogólnych problemów numerycznych w sposób znacząco wydajniejszy niż w przypadku tradycyjnych, sekwencyjnych procesorów [Cor07a]. Choć w grach komputerowych moc obliczeniową jednostek graficznych można wykorzystać do obliczeń fizyki, to CUDA idzie jeszcze dalej, umożliwiając przyspieszenie obliczeń w takich dziedzinach, jak biologia, fizyka, kryptografia, bioinformatyka oraz inne. Dla potrzeb tego segmentu NVidia opracowała specjalny procesor graficzny *Tesla* [Cor07b].

Do tej pory bardzo małe jest zainteresowanie wykorzystaniem tej technologii w procesie odkrywania wiedzy, a w szczególności znajdowania reguł asocjacyjnych. Wyniki przeprowadzonych eksperymentów pozwalają przypuszczać, że algorytm wykorzystujący możliwości procesorów wielordzeniowych będzie wyraźnie szybszy od klasycznych algorytmów eksploracji danych zaimplementowany na tradycyjnych procesorach.

1.1 Cel i zakres pracy

Celem pracy jest zaprojektowanie i zaimplementowanie algorytmu odkrywającego reguły asocjacyjne, który będzie wykorzystywał możliwości współczesnych kart graficznych dzięki wykorzystaniu technologii CUDA oraz porównanie zaprojektowanego i zaimplementowanego algorytmu do innych, podstawowych algorytmów odkrywania reguł asocjacyjnych. W ramach pracy dokonane zostanie również zebranie wiedzy dotyczącej algorytmów eksploracji reguł asocjacyjnych.

Rozdział 1 - wstęp.. Tutaj dalszy opis struktury pracy - zrobiony na koniec, gdy wszystko dalej będzie już znane.

Rozdział 2

Podstawy teoretyczne

W rozdziale tym przedstawiony zostanie przegląd literatury, który stanowi podstawy wiedzy na temat eksploracji danych, a w szczególności problemu odkrywania reguł asocjacyjnych w dużych zbiorach danych. Zebrana wiedza posłużyła autorowi do opracowania algorytmu wykorzystującego możliwości współczesnych kart graficznych.

2.1 Definicje

2.1.1 Model teoretyczny

Niech $I = \{i_1, i_2, \dots, i_m\}$ będzie zbiorem elementów o liczności $|I| = m$. *Transakcją* nazwano dowolny, niepusty podzbiór $X \subseteq I$ zbioru elementów. Bazą danych DB nazwano dowolny zbiór par (id, X) , gdzie X jest transakcją, a id jest dowolną wartością unikalną w ramach bazy danych nazywaną *identyfikatorem transakcji*. Bez utraty ogólności założono iż $id \in \mathbb{N}$.

Wsparciem (ang. *support*) $sup(X)$ transakcji X , w bazie danych nazwano częstość wystąpień transakcji w bazie danych. Formalnie przedstawia to wzór 2.1.

$$sup(X) = \frac{|\{id : (id, Y) \in DB \wedge X \subseteq Y\}|}{|DB|} \quad (2.1)$$

Łatwo zauważyć, że jeśli poziom ten jest niski, to oznacza to, że nie ma jednoznacznych dowodów na łączne występowanie elementów zbioru $Z = X \cup Y$, ponieważ zbiór Z występuje w niewielkiej liczbie transakcji.

Definicja 1 Niech będą dane dwie transakcje X i Y takie, że $X \cap Y = \emptyset$. Implikację postaci $X \Rightarrow Y$ nazwano regułą asocjacyjną.

Poziom ufności (ang. *confidence*) jest miarą zdefiniowaną dla implikacji reprezentowanej przez regułę asocjacyjną [EN05].

Definicja 2 *Poziom ufności* ($conf$) reguły asocjacyjnej $X \Rightarrow Y$ jest równy

$$conf(X \Rightarrow Y) = \frac{sup(X \cup Y)}{sup(X)} \quad (2.2)$$

Po analizie definicji 2 łatwo zauważyć, że poziom ufności może być interpretowany, jako estymacja prawdopodobieństwa $P(Y|X)$.

Reguły asocjacyjne zazwyczaj powinny spełniać pewne wymagania zdefiniowane przez użytkownika - minimalne wsparcie oraz minimalny poziom ufności, oznaczane odpowiednio *minsup* oraz *minconf*. Wyznaczają one dla aplikacji progi, jakie powinny spełniać zbiory oraz reguły, aby były brane pod uwagę w trakcie analizy.

Definicja 3 Zbiorem częstym $X \subseteq I$ nazywamy taki zbiór, który spełnia zależność $\text{sup}(X) \geq \text{minsup}$.

Generowanie reguł asocjacyjnych zazwyczaj sprowadza się do dwóch, niezależnych kroków:

1. Minimalne wsparcie jest używane do odnalezienia wszystkich zbiorów częstych w bazie danych DB .
2. Znalezione zbiory często oraz minimalny poziom ufności są używane do wygenerowania reguł asocjacyjnych.

Znalezienie wszystkich zbiorów częstych w bazie danych DB jest zadaniem wymagającym przeszukania wszystkich możliwych kombinacji bez powtórzeń ze zbioru I . Zbiór możliwych zbiorów elementów ma licznosc równą $2^n - 1$ (wszystkie zbiory, poza zbiorem pustym, który nie jest w tym wypadku zbiorem sensownym w znaczeniu poddania go analizie).

Warto zauważyć, że dla każdego zbioru częstego Y , każdy jego podzbiór X jest również zbiorem częstym [AIS93]. Korzystając z tej właściwości wsparcia możliwe jest w sposób efektywny znalezienie wszystkich zbiorów częstych w zadanej bazie danych - z tej zależności korzysta algorytm apriori opisany w rozdziale 2.2.1. Dodatkowo, wszystkie reguły zbudowane na podstawie zbioru częstego Y muszą spełniać warunek minimalnego wsparcia, ponieważ spełnia ten warunek zbiór Y , a suma zbiorów reguły jest zbiorem wyjściowym Y .

2.2 Aktualna wiedza

W rozdziale tym zebrana została oraz opracowana dotychczasowa wiedza (ang. *state-of-the-art*) na temat odkrywania reguł asocjacyjnych. Przedstawione zostaną dwa podstawowe algorytmy wykorzystywane w tym procesie: Apriori oraz FP-growth. W chwili obecnej te dwa algorytmy stanowią podstawę, na której budowane są nowe algorytmy, wykorzystujące możliwości współczesnych algorytmów (*Czy wymieniać tutaj przykłady algorytmów, które bazują na nich + refy do kilku artykułów?*).

2.2.1 Algorytm Apriori

Pierwszy algorytm odkrywający reguły asocjacyjne został przedstawiony w roku 1994 w pracy [AS94]. Niżej przedstawione zostaną szczegóły działania tego algorytmu nazwanego algorytmem *Apriori*.

Zagadnienie odkrywania reguł asocjacyjnych można podzielić na dwa etapy [AIS93]:

1. Odkrywanie zbiorów częstych, których wartość wsparcia jest wyższa od wartości *minsup*.
2. Generowanie reguł asocjacyjnych na podstawie znalezionych zbiorów częstych. Reguła $X \Rightarrow Y$ jest wynikiem działania algorytmu dla zbioru $Z = X \cup Y$, jeżeli spełnia ona nierówność $\text{conf}(X \Rightarrow Y) \geq \text{minconf}$. Ponieważ zbiór $Z = X \cup Y$ jest zbiorem częstym, to reguła spełnia również warunek przekraczania minimalnego wsparcia.

Na tym etapie możliwe jest tworzenie reguł, w których w zbiorze *poprzedników* (X z oznaczeń z definicji 1) jest wiele elementów oraz jeden w *następniku* (zbiór Y z definicji 1) [AIS93] lub dopuszczana jest możliwość wielu elementów również w następniku [AS94]. W niniejszej pracy analizowany jest sposób generowania reguł, w którym oba zbiory mogą być zbiorami wieloelementowymi.

W kolejnych podrozdziałach przedstawione zostaną etapy tworzące razem algorytmy Apriori.

k -zbiór	Zbiór zawierający k elementów.
L_k	Zbiór zawierający k -zbiory. Każdy zbiór zawarty w L_k zawiera dwa pola: i) zbiór oraz ii) wartość <i>support</i> .
C_k	Zbiór k -zbiorów kandydatów (potencjalnych zbiorów częstych). Każdy zbiór zawarty w C_k zawiera dwa pola: i) zbiór oraz ii) wartość <i>support</i> .

TABLICA 2.1: Oznaczenie skrótów w opisie algorytmu

Generowanie zbiorów częstych

W celu wyznaczenia zbiorów częstych algorytm dokonuje analizy bazy danych DB , by w kolejnych iteracjach generować rodziny coraz to liczniejszych zbiorów, będących zbiorami częstymi dla zadanej wartości *minsup*. Algorytm zaczyna od znalezienia wszystkich zbiorów jednoelementowych, które są zbiorami częstymi. W każdym kolejnym kroku generowane są zbiory częste na podstawie zbiorów wygenerowanych w kroku poprzednim. Proces ten jest kontynuowany do momentu aż nie zostaną znalezione żadne zbiory częste.

Algorytm generuje zbiory kandydatów jedynie na podstawie zbiorów częstych odkrytych w kroku poprzednim - co ważne generowanie ich odbywa się bez wielokrotnego przeglądania bazy danych transakcji. Intuicja podpowiada, że każdy podzbiór zbioru częstego jest zbiorem częstym. Zatem, każdy zbiór częsty zawierający k elementów może być wygenerowany na podstawie połączenia dwóch zbiorów posiadających $k - 1$ elementów, a na koniec kasując te zbiory, których jakikolwiek podzbiór nie jest częsty [AS94].

Tabela ?? zawiera spis oznaczeń używanych w opisie algorytmu.

Procedura APRIORI FREQUENT SET GENERATION przedstawia pseudokod realizujący opisany w tym rozdziale algorytm generowania zbiorów częstych.

APRIORI FREQUENT SET GENERATION

```

1   $L_1 \leftarrow \{1\text{-zbiory częste}\}$ 
2  for ( $k = 2; L_{k-1} \neq \emptyset; k++$ )
3      do
4           $C_k \leftarrow \text{aprioriGen}(L_{k-1})$ 
5          forall transakcja  $t \in DB$ 
6              do
7                   $C_t \leftarrow \text{subset}(C_k, t)$ 
8                  forall kandydat  $c \in C_t$ 
9                      do  $c.\text{count}++$ 
10          $L_k \leftarrow \{c \in C_k | c.\text{count} \geq \text{minsup}\}$ 
11   $\text{Answer} \leftarrow \bigcup_k L_k$ 
```

Procedura aprioriGen Procedura *aprioriGen* reprezentuje proces tworzenia zbiorów k -elementowych kandydatów na podstawie zbiorów wejściowych $k - 1$ -elementowych. Procedura ta jest podzielona na dwa etapy: łączenia oraz przycinania.

Jak łatwo zauważyć wynikiem działania JOIN STEP są zbiory k -elementowe, które powstały na podstawie zbiorów wejściowych L_{k-1} , a ich zawartość różni się tylko jednym elementem - ostatnim. Ważnym faktem jest to, iż elementy w zbiorach są uporządkowane leksykograficznie, co wykorzystywane jest w tej procedurze.

JOIN STEP

```

1 insert into  $C_k$ 
2 select  $p.item_1, p.item_2, \dots, p.item_{k-1}, q.item_{k-1}$ 
3 from  $L_{k-1} p, L_{k-1} q$ 
4 where  $p.item_1 = q.item_1, \dots, p.item_{k-2} = q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$ 

```

Warto zauważyć, że JOIN STEP jest ekwiwalentem rozszerzania zbioru L_{k-1} każdym elementem zbioru elementów I , a następnie kasowania tych $(k-1)$ -zbiorów otrzymanych przez usuwanie $(k-1)$ elementu, które nie są w L_{k-1} .

Warunek $p.item_{k-1} < q.item_{k-1}$ zapewnia, że nie będą generowane duplikaty. Dlatego też po etapie łączenia zachodzi zależność $C_k \supseteq L_k$.

Następnym krokiem jest PRUNE STEP, w którym usuwane są wszystkie elementy $c \in C_k$, którego jakiegokolwiek podzbiór $(k-1)$ -elementowy zbioru c nie należy do L_{k-1} .

PRUNE STEP

```

1 forall zbiór  $c \in C_k$ 
2   do
3     forall  $(k-1)$ -podzbiór  $s$  zbioru  $c$ 
4       do
5         if  $s \notin L_{k-1}$ 
6           then
7             delete  $c$  z  $C_k$ 

```

Celem operacji przycinania (ang. *prune*) jest ograniczenie rozmiaru zbioru C_k przed sprawdzeniem wsparcia dla kandydatów w bazie danych DB . W tym celu wykorzystywana jest właściwość, z której wynika, że jeśli jakiś $(k-1)$ -podzbiór danego kandydata ($c \in C_k$) nie występuje w L_{k-1} , to kandydat c nie jest zbiorem częstym i powinien być usunięty z C_k .

Generowanie reguł asocjacyjnych

Po zakończeniu pierwszego etapu algorytm przystępuje do drugiego, czyli do budowania reguł asocjacyjnych na podstawie odkrytych zbiorów. Podobnie, jak w [AS94] algorytm będący przedmiotem analizy niniejszej pracy, generuje wszystkie możliwe reguły asocjacyjne dla zadanego zbioru. Mniej ogólny sposób generowania reguł został przedstawiony w pracy [AIS93], jednakże podjęto decyzję, że jest to sposób zbyt mało użyteczny w środowisku produkcyjnym.

Aby wygenerować reguły, dla każdego zbioru częstego l znajdowane są niepuste podzbiory - podzbiór taki oznaczony jest jako a . Dla takich oznaczeń wygenerowana zostanie reguła $a \Rightarrow (l-a)$, jeżeli spełniona jest nierówność $\frac{support(l)}{support(a)} \geq minconr$. Warto zauważyć, że dla każdego zbioru częstego generowane są wszystkie możliwe niepuste podzbiory - zapewnia to, że odkryte zostaną wszystkie możliwe reguły.

Procedura GENERATE FREQUENT ITEMSETS prezentuje generowanie reguł asocjacyjnych na podstawie odkrytych zbiorów częstych.

GENERATE FREQUENT ITEMSETS

```

1 forall zbiór częsty  $l_k, k \geq 2$ 
2   call  $genrules(l_k, l_k)$ 

```

GENRULES(l_k : k -zbiór częsty, a_m : m -zbiór częsty)

```

1   $A \leftarrow \{(m-1)\text{-zbiór } a_{m-1} | a_{m-1} \subset a_m\}$ 
2  for  $a_{m-1} \in A$ 
3       $conf \leftarrow \frac{support(l_k)}{support(a_{m-1})}$ 
4      if  $conf \geq minconf$ 
5          output reguła  $a_{m-1} \Rightarrow (l_k - a_{m-1})$ 
          ufność =  $conf$  oraz wsparcie =  $support(l_k)$ 
6          if  $then1 > 1$ 
7              call genrules( $l_k, a_{m-1}$ )
          generowanie reguł podzbiorów zbioru  $a_{m-1}$ 
```

2.2.2 Algorytm FP-growth

Test [HPYM04]

Literatura

- [Agr94] Rakesh Agrawal. Quest: a project on database mining. *Proceedings of the ACM SIGMOD International Conference on Management of Data*, 1994.
- [AIS93] Rakesh Agrawal, Tomasz Imielinski, Arun Swami. Mining association rules between sets of items in large databases. *Proceedings of the ACM SIGMOD International Conference on Management of Data*, strony 207–216, 1993.
- [AS94] Rakesh Agrawal, Ramakrishnan Srikant. Fast algorithms for mining association rules. *VLDB*, 1994.
- [Cor07a] NVIDIA Corporation. Cuda - zone. [on-line]
http://www.nvidia.pl/object/cuda_home_new_pl.html, 2007.
- [Cor07b] NVIDIA Corporation. Nvidia tesla 20. [on-line]
http://www.nvidia.pl/object/cuda_home_new_pl.html, 2007.
- [EN05] Ramez Elmasri, Shamkand B. NAvathe. *Wprowadzenie do systemów baz danych*. Addison-Wesley, Reading, MA, USA, 2005.
- [FPSSU96] Usama M. Fayyad, Gregory Piatetsky-Shapiro, Padhraic Smyth, Ramasamy Uthurusamy. *Advances in Knowledge Discovery and Data Mining*. AAAI/MIT Press, 1996.
- [HPYM04] Jiawei Han, Jian Pei, Yiwen Yin, Runying Mao. Mining frequent patterns without candidate generation: A frequent-pattern tree approach. *Data Mining and Knowledge Discovery*, 2004.



© 2011 Tomasz Kujawa

Instytut Informatyki, Wydział Informatyki i Zarządzania
Politechnika Poznańska

Skład przy użyciu systemu L^AT_EX.

BibT_EX:

```
@mastersthesis{ key,  
  author = "Tomasz Kujawa",  
  title = "{Równoległe odkrywanie reguł asocjacyjnych zaimplementowane na procesory graficzne}",  
  school = "Poznan University of Technology",  
  address = "Pozna{\n}, Poland",  
  year = "2011",  
}
```