



UNIwersytet PRzyroDniczy we Wroclawiu

Data Mining Wykład 2

Odkrywanie asocjacji

Plan wykładu

- Wprowadzenie
- Sformułowanie problemu
- Typy reguł asocjacyjnych
- Proces odkrywania reguł asocjacyjnych

UNIwersytet PRzyroDniczy we Wroclawiu

Geneza problemu

- Geneza problemu odkrywania reguł asocjacyjnych:

problem analizy koszyka zakupów
(MBA – Market Basket Analysis)

- Dane:

baza danych zawierająca informacje o zakupach realizowanych przez klientów supermarketu

- Cel:

znalezienie grup produktów, które klienci supermarketu najczęściej kupują razem

UNIwersytet PRzyroDniczy we Wroclawiu

Analiza koszyka zakupów

- Cel analizy MBA:

znalezienie naturalnych wzorców zachowań konsumentów

- Wykorzystanie wzorców zachowań:

organizacji półek w supermarkecie

opracowania akcji promocyjnych

opracowania katalogu oferowanych produktów

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Zastosowanie MBA

- Znaleziony wzorzec:

„ktoś kto kupuje pieluski, najczęściej kupuje również piwo”

- Akcja promocyjna: (typowy trick)

Ogłoś obniżkę cen pieluszek, jednocześnie podnieś piwa

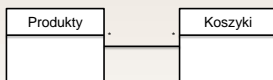
- Organizacja sklepu:

Staraj się umieszczać produkty kupowane wspólnie w przeciwległych końcach sklepu, zmuszając klientów do przejścia przez cały sklep

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Model koszyka zakupów

Model koszyka zakupów jest pewną abstrakcją umożliwiającą modelowanie relacji wiele-do-wiele pomiędzy encjami „produkty” i „koszyki”



Formalnie, model koszyka zakupów można opisać za pomocą tzw. tablicy obserwacji

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Tablica obserwacji (1)

- Dany jest zbiór atrybutów $A = \{A_1, A_2, \dots, A_n\}$ oraz zbiór obserwacji $T = \{T_1, T_2, \dots, T_m\}$

TR_{ID}	A_1	A_2	A_3	A_4	A_5
T_1	1	0	0	0	1
T_2	1	1	1	1	1
T_3	0	0	1	1	0
T_4	0	1	0	0	0
T_5	1	0	0	1	0
T_6	0	0	1	0	0
T_7	1	1	1	0	0
T_8	1	1	0	0	1

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Tablica obserwacji (2)

- Elementy tablicy obserwacji:

Atrybuty tablicy reprezentują wystąpienia encji „produkty”

Wiersze tablicy reprezentują wystąpienia encji „koszyki”

Dodatkowy atrybut TR_{ID} – wartościami atrybutu są identyfikatory poszczególnych obserwacji

Pozycja $T_i[A_j] = 1$ tablicy wskazuje, że i -ta obserwacja zawiera wystąpienie j -tego atrybutu

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Tablica obserwacji - przykłady

- „koszyki” = studenci, „produkty” = oferowane wykłady

MBA – poszukiwanie wykładów, które studenci wybierają najczęściej łącznie

- „koszyki” = strony WWW, „produkty” = słowa kluczowe

MBA – poszukiwanie stron WWW opisanych tymi samymi, lub podobnymi lub podobnymi, zbiorami słów kluczowych (prawdopodobnie, znalezione strony dotyczą podobnej problematyki)

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Reguły asocjacyjne (1)

- Wynikiem analizy koszyka jest zbiór reguł asocjacyjnych postaci następującej relacji:

$$\{(A_{i1} = 1) \wedge \dots \wedge (A_{ik} = 1)\} \rightarrow \{(A_{ik+1} = 1) \wedge \dots \wedge (A_{ik+l} = 1)\}$$

Interpretacja reguły:

„Jeżeli klient kupił produkty $A_{i1}, A_{i2}, \dots, A_{ik}$, to prawdopodobnie kupił również produkty $A_{ik+1}, A_{ik+2}, \dots, A_{ik+l}$ ”

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

Reguły asocjacyjne (2)

- Regułę asocjacyjną (1) można przedstawić jednoznacznie w równoważnej postaci:

$$\theta \rightarrow \varphi: (A_{i1}, A_{i2}, \dots, A_{ik}) \rightarrow (A_{ik+1}, A_{ik+2}, \dots, A_{ik+l})$$

- Z każdą regułą asocjacyjną $\theta \rightarrow \varphi$ związane są dwie podstawowe miary określające statystyczną ważność i siłę reguły:

$$\text{Wsparcie} - \sup(\theta \rightarrow \varphi)$$

$$\text{Ufność} - \text{conf}(\theta \rightarrow \varphi)$$

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

Reguły asocjacyjne (3)

- Statystyczna ważność i siła reguły:

Wsparciem (\sup) reguły asocjacyjnej $\theta \rightarrow \varphi$ nazywać będziemy **stosunek liczby obserwacji, które spełniają warunek $\theta \wedge \varphi$, do liczby wszystkich obserwacji** (wsparcie reguły = prawdopodobieństwu zajścia zdarzenia $\theta \wedge \varphi$)

Ufnością (conf) reguły asocjacyjnej $\theta \rightarrow \varphi$ nazywać będziemy **stosunek liczby obserwacji, które spełniają warunek $\theta \wedge \varphi$, do liczby obserwacji, które spełniają warunek θ** (ufność reguły = warunkowemu prawdopodobieństwu $p(\varphi | \theta)$)

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

Klasyfikacja reguł asocjacyjnych

- Klasyfikacja reguł asocjacyjnych ze względu na:

Typ przetwarzanych danych

Wymiary przetwarzanych danych

Stopień abstrakcji przetwarzanych danych

- Inne typy reguł asocjacyjnych
- Asocjacje vs. analiza korelacji

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Typ przetwarzanych danych (1)

- Wyróżniamy:

binarne reguły asocjacyjne - regułę asocjacyjną nazywamy binarną, jeżeli dane występujące w regule są danymi (zmiennymi) binarnymi

ilościowe reguły asocjacyjne - regułę asocjacyjną nazywamy ilościową, jeżeli dane występujące w regule są danymi ciągłymi i/lub kategorycznymi

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Typ przetwarzanych danych (2)

- Binarna reguła asocjacyjna:

pieluszki = 1 \rightarrow piwo = 1

- reprezentuje współwystępowanie danych

- Ilościowa reguła asocjacyjna:

wiek = '30...40' \wedge wykształcenie = 'wyższe' \rightarrow
opcja_polityczna = 'demokrata'

- reprezentuje współwystępowanie wartości danych

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Wymiarowość przetwarzanych danych (1)

- Wyróżniamy:

jednowymiarowe reguły asocjacyjne - regułę asocjacyjną nazywamy jednowymiarową, jeżeli dane występujące w regule reprezentują tę samą dziedzinę wartości.

wielowymiarowe reguły asocjacyjne - regułę asocjacyjną nazywamy wielowymiarową, jeżeli dane występujące w regule reprezentują różne dziedziny wartości.

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

Typ przetwarzanych danych (2)

- Jednowymiarowa reguła asocjacyjna:

pieluszki = 1 → piwo = 1

- Wielowymiarowa reguła asocjacyjna:

wiek = '30...40' ∧ wykształcenie = 'wyższe' →
opcja_polityczna = 'demokrata'

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

Stopień abstrakcji przetwarzanych danych (1)

- Wyróżniamy:

jednopoziomowe reguły asocjacyjne - regułę asocjacyjną nazywamy jednopoziomą, jeżeli dane występujące w regule reprezentują ten sam poziom abstrakcji.

Wielopoziomowe reguły asocjacyjne - regułę asocjacyjną nazywamy wielopoziomą, jeżeli dane występujące w regule reprezentują różne poziomy abstrakcji.

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

Stopień abstrakcji przetwarzanych danych (2)

- Jednopoziomowa reguła asocjacyjna:

$\text{pieluszki_Pampers} = 1 \rightarrow \text{piwo_Zywiec} = 1$

- Wielopoziomowa reguła asocjacyjna:

$\text{pieluszki_Pampers} = 1 \wedge \text{piwo_Zywiec} = 1 \rightarrow \text{napoje} = 1$

(produkt napoje reprezentuje pewną abstrakcję, będącą generalizacją określonych produktów)

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Odkrywanie binarnych reguł asocjacyjnych

- Dane:
 - $I = \{i_1, i_2, \dots, i_n\}$: zbiór literałów, nazywanych dalej elementami
 - Transakcja T : zbiór elementów, takich że $T \subseteq I$ i $T \neq \emptyset$
 - Baza danych D : zbiór transakcji
- Transakcja T wspiera element $x \in I$, jeżeli $x \in T$
- Transakcja T wspiera zbiór $X \subseteq I$, jeżeli T wspiera każdy element ze zbioru X , $X \subseteq T$

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Reguły asocjacyjne – miary (1)

- Binarna reguła asocjacyjna:

Binarną regułą asocjacyjną (krótko, regułą asocjacyjną) nazywamy relację postaci $X \rightarrow Y$, gdzie $X \subseteq I$, $Y \subseteq I$, i $X \cap Y = \emptyset$

- Wsparcie (support):

Reguła $X \rightarrow Y$ posiada wsparcie sup w bazie danych D , $0 \leq \text{sup} \leq 1$, jeżeli $\text{sup}\%$ transakcji w D wspiera zbiór $X \cup Y$

- Ufność (confidence):

Reguła $X \rightarrow Y$ posiada ufność conf w bazie danych D , $0 \leq \text{conf} \leq 1$, jeżeli $\text{conf}\%$ transakcji w D , które wspierają zbiór X , wspierają również Y

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Reguły asocjacyjne – miary (2)

- wsparcie($X \rightarrow Y$):

oznacza liczbę transakcji w bazie danych, które potwierdzają daną regułę – miara wsparcia jest symetryczna względem zbiorów stanowiących poprzednik i następnik reguły

- ufność($X \rightarrow Y$):

oznacza stosunek liczby transakcji zawierających $X \cup Y$ do liczby transakcji zawierających Y – miara ta jest asymetryczna względem zbiorów stanowiących poprzednik i następnik reguły

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Reguły asocjacyjne – miary (3)

- Ograniczenia miar (definiowane przez użytkownika):

Minimalne wsparcie – minsup

Minimalna ufność – minconf

- Mówimy, że reguła asocjacyjna $X \rightarrow Y$ jest silna jeżeli:

$\text{sup}(X \rightarrow Y) \geq \text{minsup}$ i $\text{conf}(X \rightarrow Y) \geq \text{minconf}$

- Dana jest baza danych transakcji Należy znaleźć wszystkie silne binarne reguły asocjacyjne

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Przykład

TR _{id}	Produkty
1	A,B,C
2	A,C
3	A,D
4	B,E,F

Zakładając:

minsup = 50% oraz minconf = 50%

w przedstawionej bazie danych można znaleźć następujące reguły asocjacyjne:

$A \rightarrow C$ sup = 50%, conf = 66,6 %
 $C \rightarrow A$ sup = 50%, conf = 100%

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

Podsumowanie

- Typy reguł asocjacyjnych
 - Typ przetwarzanych danych
 - Wymiarowość przetwarzanych danych
 - Stopień abstrakcji przetwarzanych danych
- Proces odkrywania reguł asocjacyjnych
 - Wsparcie - $\text{sup}(\theta \rightarrow \varphi)$
 - Ufność - $\text{conf}(\theta \rightarrow \varphi)$
