

# Eksploracja danych

reguły asocjacyjne

Piotr Lipiński

## Reguły asocjacyjne

- Badania dotyczące reguł asocjacyjnych zostały rozpoczęte na potrzeby analizy transakcji dokonywanych przez klientów w supermarketach (tzw. analiza koszyka zakupów).
- Rozpatrujemy pewien zbiór przedmiotów  $\mathbf{I} = \{I_1, I_2, \dots, I_d\}$  i pewien zbiór transakcji  $\mathbf{T} = \{T_1, T_2, \dots, T_N\}$ , gdzie każda transakcja  $T_i$  to pewien podzbiór zbioru przedmiotów  $\mathbf{I}$ .
- Co jest interesujące ?
  - jak często były kupowane poszczególne przedmioty
    - dokładniej: w ilu procentach transakcji występowały poszczególne przedmioty
      - odpowiedzi dostarcza prosta statystyka
  - jakie przedmioty były kupowane razem
    - dokładniej: jakie są częste zbiory przedmiotów występujące w transakcjach
      - odpowiedź wymaga nietrywialnej analizy danych
  - kupno jakiego przedmiotu pociąga za sobą kupno przedmiotu  $X$

## Reguły asocjacyjne

□ Przykład:

1	plums, lettuce, tomatoes
2	celery, confectionery
3	confectionery
4	apples, carrots, tomatoes, potatoes, confectionery
5	apples, oranges, lettuce, tomatoes, confectionery
6	peaches, oranges, celery, potatoes
7	beans, lettuce, tomatoes
8	oranges, lettuce, carrots, tomatoes, confectionery
9	apples, bananas, plums, carrots, tomatoes, onions, confectionery
10	apples, potatoes

Piotr Lipiński, Eksploracja danych

3

## Reguły asocjacyjne

□ Dla zbioru przedmiotów  $X \subseteq I$  określamy

- support  $X$  (nośnik, wsparcie)

$$\text{supp}(X) = |\{T \in \mathbf{T} : X \subseteq T\}| / N.$$

□ Częste zbiory przedmiotów

- Zbiór przedmiotów  $X \subseteq I$  nazywamy częstym, jeśli występuje w dużej liczbie transakcji, czyli  $\text{supp}(X) > \alpha$ , dla pewnego ustalonego progu  $\alpha > 0$ .
- Fakt 1: Jeśli zbiór przedmiotów  $X \subseteq I$  jest częsty, to każdy jego podzbiór  $Y \subseteq X$  też jest częsty.
- Fakt 2: Jeśli zbiór przedmiotów  $X \subseteq I$  zawiera podzbiór  $Y \subseteq X$ , który nie jest częsty, to sam zbiór  $X$  też nie jest częsty.
- W związku z tym konstrukcję wszystkich zbiorów częstych można rozpocząć od zbiorów częstych jednoelementowych, a następnie łączyć je kolejno w zbiory częste dwuelementowe, trzelementowe, itd.

Piotr Lipiński, Eksploracja danych

4

## Reguły asocjacyjne

- Reguła asocjacyjna  $A \rightarrow C$  to para zbiorów przedmiotów  $A \subseteq I$  i  $C \subseteq I$ .
  - Interpretacja: Jeśli w transakcji  $T$  znalazły się przedmioty ze zbioru  $A$ , tzn.  $A \subseteq T$ , to z dużym prawdopodobieństwem w transakcji  $T$  znalazły się przedmioty ze zbioru  $B$ , tzn.  $B \subseteq T$ .
- Dla reguły asocjacyjnej  $A \rightarrow C$  określamy
  - support  $A \rightarrow C$  (nośnik, wsparcie)
$$\text{supp}(A \rightarrow C) = \text{supp}(A \cup C),$$
  - confidence  $A \rightarrow C$  (wiarygodność, zaufanie)
$$\text{conf}(A \rightarrow C) = \text{supp}(A \cup C) / \text{supp}(A).$$
- Jeśli  $\text{supp}(A \rightarrow C)$  jest wystarczająco wysokie, a zarejestrowane transakcje reprezentują próbę losową z tego samego rozkładu co przyszłe transakcje, to  $\text{conf}(A \rightarrow C)$  jest dość dobrym estymatorem prawdopodobieństwa, że dowolna przyszła transakcja zawierająca  $A$  będzie zawierać także  $B$ .
- Celem jest wyznaczenie zbioru wszystkich reguł asocjacyjnych o wystarczająco wysokim wsparciu i jak najwyższej wiarygodności.

Piotr Lipiński, Eksploracja danych

5

## Algorytm apriori

- IDEA ALGORYTMU:
  - jako parametr algorytmu określamy minimalne wsparcie min-supp,
  - zbiory przedmiotów i reguły asocjacyjne o wsparciu poniżej minimalnego uznajemy za nieinteresujące (niewiarygodne),
  - skoro  $\text{supp}(A \rightarrow C) = \text{supp}(A \cup C)$ , to interesujące są tylko takie reguły asocjacyjne  $A \rightarrow C$ , dla których  $A \cup C$  jest zbiorem częstym (o minimalnym wsparciu min-supp),
  - generujemy wszystkie zbiory częste  $X$ ,
  - dla każdego zbioru częstego  $X$ , rozpatrujemy wszystkie reguły asocjacyjne  $A \rightarrow C$ , dla których  $A \cup C = X$ , tzn. generujemy rozbicia zbioru  $X$  na podzbiory  $A$  i  $C$ ,
  - ograniczamy się do rozbicia zbioru  $X$  na ROZŁĄCZNE podzbiory  $A$  i  $C$ .

Piotr Lipiński, Eksploracja danych

6

## Algorytm apriori

### □ ALGORYTM A PRIORI:

```
k = 1; Lk = zbiór jednoelementowych zbiorów częstych;
k = 2;
while Lk-1 ≠ ∅ do
    Ck = { {x1, x2, ..., xk-2, xk-1, xk} :
            {x1, x2, ..., xk-2, xk-1} ∈ Lk-1 ∧
            {x1, x2, ..., xk-2, xk} ∈ Lk-1 };
    for each t ∈ T do
        for each c ∈ Ck ∧ c ⊆ t do
            c.count++;
    Lk = {c ∈ Ck : c.count ≥ min-supp};
return L1 ∪ L2 ∪ ... ∪ Lk;
```

Piotr Lipiński, Eksploracja danych

7

## Przykład

- Niech min-confidence = 0.75.
- Niech X = {lettuce, tomatoes}.
- Łatwo policzyć, że support(X) = 0.4.
- S<sub>0</sub> = { {lettuce, tomatoes} → ∅ }.
- Rozpatrujemy:
  - confidence({lettuce} → {tomatoes}) = 0.4 / 0.4
  - confidence({tomatoes} → {lettuce}) = 0.4 / 0.6
- S<sub>1</sub> = { {lettuce} → {tomatoes} }
- Rozpatrujemy:
  - confidence(∅ → {lettuce, tomatoes}) = 0.4
- S<sub>2</sub> = ∅
- Ostatecznie, wiarygodne reguły uzyskane ze zbioru częstego X, to:
  - {lettuce, tomatoes} → ∅
  - {lettuce} → {tomatoes}

Piotr Lipiński, Eksploracja danych

8

## Algorytm apriori

•Generowanie reguł asocjacyjnych ze zbioru częstego  $X$  oparte jest na rozbijaniu tego zbioru w różny sposób na dwa rozłączne podzbiory  $Y$  i  $Z$ , takie że  $X = Y \cup Z$ , i rozpatrywaniu reguł  $Y \rightarrow Z$ . Zwracane są tylko wiarygodne reguły  $Y \rightarrow Z$ , czyli takie których  $\text{confidence}(Y \rightarrow Z) \geq \text{min\_confidence}$ .

•**Spostrzeżenie 1:**  $\text{confidence}(A \rightarrow B \cup C) \leq \text{confidence}(A \cup B \rightarrow C)$

•Dowód: wprost z definicji confidence, bo  $\text{support}(A) \geq \text{support}(A \cup B)$ .

•**Wniosek:** Jeśli reguła  $A \rightarrow B \cup C$  jest wiarygodna, to reguła  $A \cup B \rightarrow C$  tym bardziej jest wiarygodna.

•**Spostrzeżenie 2:**  $\text{confidence}(X \rightarrow \emptyset) = 1$  i żadna inna reguła powstała z  $X$  nie ma większego confidence.

•**Spostrzeżenie 3:**  $\text{confidence}(\emptyset \rightarrow X) = \text{support}(X)$  i żadna inna reguła powstała z  $X$  nie ma mniejszego confidence.

•**Konstruowanie reguł:**

- Zaczynamy od reguły  $X \rightarrow \emptyset$ . Niech  $k = 0$ , zaś zbiór reguł wiarygodnych  $S_k = \{X \rightarrow \emptyset\}$ .
- Będziemy próbować przesuwac przedmioty z lewej strony na prawą. Niech  $S_{k+1} = \emptyset$ . Dla każdej reguły  $A \rightarrow B$  z  $S_k$ , sprawdzamy wszystkie reguły postaci  $A \setminus \{x\} \rightarrow B \cup \{x\}$  (dla każdego  $x \in A$ ). Jeżeli reguła jest wiarygodna, to dodajemy ją do  $S_{k+1}$ . W tym kroku sprawdzamy więc  $|S_k|(|X| - k)$  reguł.
- Niech  $k = k + 1$ . Jeśli  $S_k$  jest niepusty, to wracamy do kroku 1.
- Zwracamy wszystkie reguły ze zbiorów  $S_0, S_1, \dots, S_k$ .

## Które reguły są interesujące?

- Często otrzymujemy wiele reguł, z których nie wszystkie są interesujące.
- Często bada się różnicę między wsparciem reguły, a iloczynem wsparcia lewej i prawej strony reguły.
- Jeśli wsparcie reguły jest (w przybliżeniu) równe iloczynowi wsparcia lewej i prawej strony, to mówimy, że lewa i prawa strona są niezależne. Takie reguły są nieinteresujące, niezależnie od ich wiarygodności.
  - Dlaczego?
- Powstała więc idea stworzenia uniwersalnych miar oceniających jakość reguł asocjacyjnych.

## Które reguły są interesujące?

- Dla reguły asocjacyjnej  $A \rightarrow C$  określamy dodatkowe miary

- lift  $A \rightarrow C$

$$\text{lift}(A \rightarrow C) = \text{confidence}(A \rightarrow C) / \text{supp}(C)$$

- leverage  $A \rightarrow C$

$$\text{leverage}(A \rightarrow C) = \text{supp}(A \rightarrow C) - \text{supp}(A) \text{supp}(C).$$

## Które reguły są interesujące?

- Przykład:

$\{\text{pomidory}\} \rightarrow \{\text{sałata}\}$

$\text{support}(\{\text{sałata}\}) = 0.4$

$\text{confidence}(\{\text{pomidory}\} \rightarrow \{\text{sałata}\}) = 0.67$

$\text{lift}(\{\text{pomidory}\} \rightarrow \{\text{sałata}\}) = 0.67 / 0.4 = 1.675$

$\{\text{pomidory}\} \rightarrow \{\text{słodycze}\}$

$\text{support}(\{\text{słodycze}\}) = 0.6$

$\text{confidence}(\{\text{pomidory}\} \rightarrow \{\text{słodycze}\}) = 0.67$

$\text{lift}(\{\text{pomidory}\} \rightarrow \{\text{słodycze}\}) = 0.67 / 0.6 = 1.117$

WNIOSEK: Pomidory mają większy wpływ na sałatę niż na słodycze.

## Które reguły są interesujące?

□ Przykład:

$\{\text{marchew}\} \rightarrow \{\text{pomidory}\}$ , confidence = 1, lift = 1.667

$\{\text{sałata}\} \rightarrow \{\text{pomidory}\}$ , confidence = 1, lift = 1.667

*UWAGA: Jednak druga reguła może być bardziej interesująca, bo dotyczy większej liczby transakcji.*

$\text{support}(\{\text{marchew}\} \rightarrow \{\text{pomidory}\}) = 0.3$

$\text{support}(\{\text{marchew}\}) = 0.3$

$\text{support}(\{\text{pomidory}\}) = 0.6$

$\text{leverage}(\{\text{marchew}\} \rightarrow \{\text{pomidory}\}) = 0.3 - 0.3 * 0.6 = 0.12$

$\text{support}(\{\text{sałata}\} \rightarrow \{\text{pomidory}\}) = 0.4$

$\text{support}(\{\text{sałata}\}) = 0.4$

$\text{leverage}(\{\text{sałata}\} \rightarrow \{\text{pomidory}\}) = 0.4 - 0.4 * 0.6 = 0.16$

## Reguły podobne

□ Zobaczmy dwie wykryte reguły

$\text{apples} \rightarrow \text{confectionery, tomatoes}$

$\text{apples, tomatoes} \rightarrow \text{confectionery}$

obie one mają takie samo wsparcie, równe 0.3, ale różną wiarygodność.

□ Obie reguły są partycją tego samego częstego zbioru przedmiotów

$\{\text{apples, confectionery, tomatoes}\}$ ,

więc np. przy ustawianiu produktów na sklepowych półkach przydatność obu reguł jest taka sama.

□ Jak oceniać jakość wykrytych zbiorów przedmiotów?

## Efektywność algorytmu

- Algorytm APRIORI jest dosyć efektywny w przypadku danych rzadkich (sparse data), kiedy każdy przedmiot występuje w niewielkiej liczbie transakcji.
- Istnieją modyfikacje tego algorytmu dla innych sytuacji (praca Bayardo 1997), opierają się na różnych sposobach konstrukcji częstych zbiorów przedmiotów.
- Istnieją też modyfikacje algorytmu konstruujące reguły asocjacyjne bez uprzedniego konstruowania częstych zbiorów przedmiotów.

## Algorytm PrefixSpan

- Dotychczas interesowały nas transakcje będące zbiorami przedmiotów. Tak postrzegane transakcje, jako zbiory, nie zawierały informacji o kolejności występowania przedmiotów.
- Interesującym rozszerzeniem tradycyjnego podejścia jest określenie transakcji jako ciągów przedmiotów, co umożliwi określenie kolejności występowania przedmiotów.
- Potencjalne zastosowanie: śledzenie aktywności użytkowników w portalu internetowym (przedmiot = strona portalu, transakcja = ciąg odwiedzonych stron portalu).



## Ciągi częste

□ Przykład:

$I = \{a, b, c, d\}$

$D = \{ \langle 1, \langle abc \rangle \rangle, \langle 2, \langle ac \rangle \rangle, \langle 3, \langle bd \rangle \rangle, \langle 4, \langle acd \rangle \rangle \}$

$\text{min-supp} = 3$

Ciągi częste to  $\{\langle a \rangle, \langle c \rangle, \langle ac \rangle\}$

□ Interpretacja:

- Portal zawiera cztery strony internetowe a, b, c, d.
- Zarejestrowano następujące wizyty użytkowników: abc, ac, bd, acd.
- Wówczas często odwiedzane strony portalu to a i c, a często wykonywana ścieżka w portalu to przejście ze strony a (niekoniecznie bezpośrednio) do strony c.

## Ciągi częste

- Zamiast częstych zbiorów i budowanych na nich reguł asocjacyjnych zajmijmy się przypadkiem częstych ciągów (w których ważna jest kolejność elementów)
- Czym się różni to podejście od częstych zbiorów?
- Czy można zastosować algorytm a priori do konstrukcji ciągów częstych?
- Jednym z algorytmów do wyszukiwania ciągów częstych jest algorytm PrefixSpan.

## Ciągi częste

- $I = \{a, b, c, d, e\}$
- $D = \{ \langle 1, \langle abc \rangle \rangle, \langle 2, \langle acbed \rangle \rangle, \langle 3, \langle bd \rangle \rangle, \langle 4, \langle bcdad \rangle \rangle, \langle 5, \langle becabc \rangle \rangle \}$
- $\text{min-supp} = 3$
  
- Krok 1: Znajdź zbiór 1-elementowych ciągów częstych.
  - Wystarczy jednokrotnie przejrzeć zbiór transakcji  $D$  i policzyć liczby wystąpień poszczególnych przedmiotów z  $I$ , a następnie wybrać te przedmioty, które występują co najmniej  $\text{min-supp}$  razy.
  - $L_1 = \{ \langle a \rangle, \langle b \rangle, \langle c \rangle, \langle d \rangle \}$
- Krok 2: Podziel przestrzeń obliczeń.
  - Dłuższe ciągi częste muszą rozpoczynać się prefiksem będącym jednym ze znalezionych 1-elementowych ciągów częstych.
  - Szukamy ciągów częstych o prefiksie  $\langle a \rangle$ , prefiksie  $\langle b \rangle$ , prefiksie  $\langle c \rangle$  i prefiksie  $\langle d \rangle$ .
- Krok 3: Znajdź kolejne zbiory ciągów częstych.
  - Kolejne ciągi częste można wyznaczyć rzutując zbiór transakcji  $D$  na poszczególne 1-elementowe prefiksy z  $L_1$ .
  - $D|_{\langle a \rangle} = \{ \langle 1, \langle bc \rangle \rangle, \langle 2, \langle cbcd \rangle \rangle, \langle 4, \langle d \rangle \rangle, \langle 5, \langle bc \rangle \rangle \}$

## Ciągi częste

- Krok 3a: Znajdź zbiór ciągów częstych o prefiksie  $\langle a \rangle$ 
  - Rzutujemy zbiór transakcji  $D$  na prefiks  $\langle a \rangle$  uzyskując:  
 $D|_{\langle a \rangle} = \{ \langle 1, \langle bc \rangle \rangle, \langle 2, \langle cbcd \rangle \rangle, \langle 4, \langle d \rangle \rangle, \langle 5, \langle bc \rangle \rangle \}.$
  - Wyznaczamy 1-elementowe ciągi częste w  $D|_{\langle a \rangle}$  uzyskując:  
 $\langle b \rangle$  i  $\langle c \rangle.$
  - Zatem wszystkie 2-elementowe ciągi częste w  $D$  o prefiksie  $\langle a \rangle$  to:  
 $\langle ab \rangle$  i  $\langle ac \rangle.$
  - Dłuższe ciągi częste w  $D$  o prefiksie  $\langle a \rangle$  muszą mieć więc prefiks  $\langle ab \rangle$  lub  $\langle ac \rangle$ , znajdujemy je rekurencyjnie

## Ciągi częste

- Krok 3ab: Znajdź zbiór ciągów częstych o prefiksie <ab>
  - Rzutujemy zbiór transakcji D na prefiks <ab> uzyskując:
$$D|_{\langle ab \rangle} = \{ \langle 1, \langle c \rangle \rangle, \langle 2, \langle ed \rangle \rangle, \langle 5, \langle c \rangle \rangle \}.$$
  - Wyznaczamy 1-elementowe ciągi częste w  $D|_{\langle ab \rangle}$  uzyskując zbiór pusty (nie ma żadnych ciągów częstych w  $D|_{\langle ab \rangle}$ ).
  - Zatem nie ma żadnych 3-elementowych ciągów częste w D o prefiksie <ab>.
- Krok 3ac: Znajdź zbiór ciągów częstych o prefiksie <ac>
  - Rzutujemy zbiór transakcji D na prefiks <ab> uzyskując:
$$D|_{\langle ac \rangle} = \{ \langle 1, \langle e \rangle \rangle, \langle 2, \langle bed \rangle \rangle, \langle 5, \langle e \rangle \rangle \}.$$
  - Wyznaczamy 1-elementowe ciągi częste w  $D|_{\langle ac \rangle}$  uzyskując zbiór pusty (nie ma żadnych ciągów częstych w  $D|_{\langle ac \rangle}$ ).
  - Zatem nie ma żadnych 3-elementowych ciągów częste w D o prefiksie <ac>.
- Zatem wszystkie ciągi częste w D o prefiksie <a> to: <a>, <ab> i <ac>.

## Ciągi częste

- Analogicznie wyznaczamy wszystkie ciągi częste w D o prefiksie <b>, <c> i <d>.
- Ostatecznie otrzymujemy:
  - <a>, <ab>, <ac>,
  - <b>, <bc>, <bd>,
  - <c>,
  - <d>.

## PrefixSpan

### □ Algorytm PrefixSpan

```
PrefixSpan(p, D|p)
  F = {1-elementowe ciągi częste z D|p}
  R = ∅
  for each f ∈ F do
    q = p · f
    R = R ∪ {q}
    wyznacz D|q
    R = R ∪ PrefixSpan(q, D|q)
  return R
```