



Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego  
podziału

Dokonanie podziału  
SPRINT - wersja  
równoległa

# WYKŁAD 7: DRZEWA DECYZYJNE DLA DUŻYCH ZBIORÓW DANYCH

Nguyen Hung Son



## FUNKCJA REKURENCYJNA *buduj\_drzewo*( $U, dec, \mathbf{T}$ ):

- 1: **if** (*kryterium\_stopu*( $U, dec$ ) = **true**) **then**
- 2:      $\mathbf{T}.etykieta = \textit{kategoria}(U, dec)$ ;
- 3:     **return**;
- 4: **end if**
- 5:  $t := \textit{wybierz\_test}(U, \mathbf{TEST})$ ;
- 6:  $\mathbf{T}.test := t$ ;
- 7: **for**  $v \in R_t$  **do**
- 8:      $U_v := \{x \in U : t(x) = v\}$ ;
- 9:     utwórz nowe poddrzewo  $\mathbf{T}'$ ;
- 10:      $\mathbf{T}.ga\acute{z}\acute{a}z(v) = \mathbf{T}'$ ;
- 11:     *buduj\_drzewo*( $U_v, dec, \mathbf{T}'$ )
- 12: **end for**



- **Kryterium stopu:** Zatrzymamy konstrukcji drzewa, gdy aktualny zbiór obiektów:
  - jest pusty lub
  - zawiera obiekty wyłącznie jednej klasy decyzyjnej lub
  - nie ulega podziału przez żaden test
- **Wyznaczenie etykiety zasadą większościową:**

$$\text{kategoria}(P, dec) = \arg \max_{c \in V_{dec}} |P_{[dec=c]}|$$

tzn., etykietą dla danego zbioru obiektów jest klasa decyzyjna najliczniej reprezentowana w tym zbiorze.

- **Kryterium wyboru testu:** heurystyczna funkcja oceniająca testy.



# MIARY RÓŻNORODNOŚCI ZBIORU

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

Każdy zbiór obiektów  $X$  ulega podziału na klasy decyzyjne:

$$X = C_1 \cup C_2 \cup \dots \cup C_d$$

gdzie  $C_i = \{u \in X : dec(u) = i\}$ .

Wektor  $(p_1, \dots, p_r)$ , gdzie  $p_i = \frac{|C_i|}{|X|}$ , nazywamy **rozkładem klas decyzyjnych** w  $X$ .

$$Conflict(X) = \sum_{i < j} |C_i| \times |C_j| = \frac{1}{2} \left( |X|^2 - \sum |C_i|^2 \right)$$

$$Entropy(X) = - \sum \frac{|C_i|}{|X|} \cdot \log \frac{|C_i|}{|X|} = - \sum p_i \log p_i$$

$$Gini(X) = 1 - \sum p_i^2$$



# OCENA FUNKCJI TESTU

Data mining

Nguyen Hung Son

Motywacje

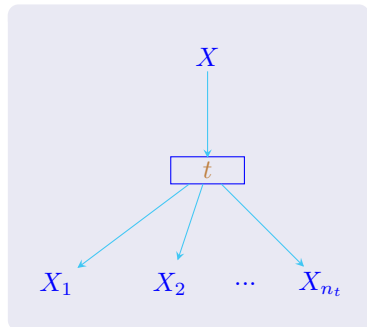
Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału

SPRINT - wersja równoległa

Każdy test  $t$  jest oceniony na podstawie informacji zawartych w  $X, X_1, \dots, X_{n_t}$



Podział zbioru  $X$  dokonany przez test  $t$ ;



Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

- Rozróżnialność:

$$disc(t, X) = conflict(X) - \sum conflict(X_i)$$

- Przyrostu informacji (Information gain).

$$Gain(t, X) = Entropy(X) - \sum \frac{|X_i|}{|X|} \cdot Entropy(X_i)$$

- Gini's index

$$G(t, X) = Gini(X) - \sum \frac{|X_i|}{|X|} \cdot Gini(X_i)$$

- Kara za zbyt drobny podział, np. gain ratio

$$Gain\_ratio = \frac{Gain(t, X)}{-\sum_{i=1}^r \frac{|X_i|}{|X|} \cdot \log \frac{|X_i|}{|X|}}$$



# PRZYKŁAD

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

Chest	No	No	No	Yes	Yes	Yes	No	No	No	No	
Taxable Income											
Sorted Values	60	70	75	85	90	95	100	120	125	220	
Split Positions	55	65	72	80	87	92	97	110	122	172	230
	<= >	<= >	<= >	<= >	<= >	<= >	<= >	<= >	<= >	<= >	<= >
Yes	0 3	0 3	0 3	0 3	1 2	2 1	3 0	3 0	3 0	3 0	3 0
No	0 7	1 6	2 5	3 4	3 4	3 4	3 4	4 3	5 2	6 1	7 0
Gini	0.420	0.400	0.375	0.343	0.417	0.400	<u>0.300</u>	0.343	0.375	0.400	0.420



# SŁABOŚCI STANDARDOWEGO ALGORYTMU:

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

- Każdy węzeł jest skojarzony z podzbiorem danych: **ograniczenie pamięciowe**
- Wyznaczenie najlepszego podziału wymaga wielokrotnego sortowania danych: **czasochłonne**
  - Dany atrybut rzeczywisty  $a$  i zbiór możliwych cięć  $(c_1, c_2, \dots, c_N)$ , najlepszy test  $(a, c_i)$  można znaleźć w czasie  $\Omega(N)$
  - Minimalna liczba prostych zapytań SQL potrzebna do szukania najlepszego testu jest  $O(dN)$ , gdzie  $d$  jest liczbą klas decyzyjnych
- Wniosek: szukanie najlepszego podziału jest kosztowne, jeśli atrybut zawiera dużo różnych wartości.





# CHARACTERYSTYKA ALGORYTMU SPRINT

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

- Nadaje się dla danych częściowo umieszczonych na dysku
- Używa się techniki pre-sortowania w celu przyspieszenia procesu obliczenia na atrybutach rzeczywistych;
- Dane są sortowane tylko raz przed obliczeniem
- Łatwo można zrównoleglić



# STRUKTURA DANYCH W SPRINT

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

- Każdy atrybut ma swoją listę wartości
- Każdy element listy ma trzy pole:
  - *wartość atrybutu*,
  - *numer klasy* i
  - *rid* (numer obiektu w zbiorze danych)
- Rzeczywiste atrybuty są uporządkowane (tylko raz przy utworzeniu)
- Na początku listy są stowarzyszone z korzeniem drzewa
- Kiedy węzeł podlega podziału, listy są podzielone i są skojarzone z odpowiednimi następnikami
- Listy są zapisane na dysku w razie potrzeby



# PRZYKŁAD

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

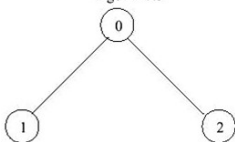
Dokonanie podziału  
SPRINT - wersja równoległa

Attribute lists for node 0

Age	Class	Tid
17	High	1
20	High	5
23	High	0
32	Low	4
43	High	2
68	Low	3

Car Type	Class	Tid
family	High	0
sports	High	1
sports	High	2
family	Low	3
truck	Low	4
family	High	5

Age < 27.5



Attribute lists for node 1

Age	Class	Tid
17	High	1
20	High	5
23	High	0

Car Type	Class	Tid
family	High	0
sports	High	1
family	High	5

Attribute lists for node 2

Age	Class	Tid
32	Low	4
43	High	2
68	Low	3

Car Type	Class	Tid
sports	High	2
family	Low	3
truck	Low	4



# STRUKTURA DANYCH W SPRINT

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

- SPRINT używa:
  - indeksu Gini do oceny jakości podziału
  - testu typu  $(a \leq c)$  dla atrybutów rzeczywistych
  - testu typu  $(a \in V)$  dla atrybutów symbolicznych
- Histogram: rozkład klas decyzyjnych zbadanego zbioru danych
- Dla atrybutu rzeczywistego dwa histogramy:
  - $C_{below}$  : histogram dla danych poniżej wartości progowej
  - $C_{above}$  : histogram dla danych powyżej wartości progowej
- Dla atrybutu symbolicznego jeden histogram zwany count matrix

# PRZYKŁAD

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

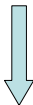
Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

<i>Car Type</i>	<i>Class</i>	<i>rid</i>
family	High	0
sports	High	1
sports	High	2
family	Low	3
truck	Low	4
family	high	5

← Punkt podziału

<i>Age</i>	<i>Class</i>	<i>rid</i>
17	High	1
20	High	5
23	High	0
32	Low	4
43	High	2
68	Low	3



Count matrix

	H	L
family	2	1
sports	2	0
truck	0	1



Histogram klas

	H	L
$C_{\text{below}}$	3	0
$C_{\text{above}}$	1	2



Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

## 1 Motywacje

## 2 ALGORYTM SPRINT

- Szukanie najlepszego podziału
- Dokonanie podziału
- SPRINT - wersja równoległa



# WYZNACZANIE PODZIAŁU ATRYBUTU RZECZYWISTEGO

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału SPRINT - wersja równoległa

Attribute List			Position of cursor in scan
Age	Class	tid	
17	High	1	← position 0
20	High	5	
23	High	0	← position 3
32	Low	4	
43	High	2	
68	Low	3	← position 6

## State of Class Histograms

		H	L
cursor position 0:	$C_{\text{below}}$	0	0
	$C_{\text{above}}$	4	2
		H	L
cursor position 3:	$C_{\text{below}}$	3	0
	$C_{\text{above}}$	1	2
		H	L
cursor position 6:	$C_{\text{below}}$	4	2
	$C_{\text{above}}$	0	0



# WYZNACZANIE PODZIAŁU ATRYBUTU SYMBOLICZNEGO

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

Lista wartości **CarType**

Car Type	Class	rid
family	High	0
sports	High	1
sports	High	2
family	Low	3
truck	Low	4
family	high	5



Count Matrix

	H	L
family	2	1
sports	2	0
truck	0	1

1. Wyznacz macierz rozkładu klas obiektów w danym węźle
2. Używając algorytmu aproksymacyjnego (w SLIQ) wyznacz podzbiór wartości  $V \subseteq D_a$  t. żeby test  $(a \in V)$  był optymalny





# OUTLINE

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

**Dokonanie podziału**

SPRINT - wersja równoległa

## 1 Motywacje

## 2 ALGORYTM SPRINT

- Szukanie najlepszego podziału
- **Dokonanie podziału**
- SPRINT - wersja równoległa



# GŁÓWNA IDEA

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

- Każda lista jest podzielona na dwie listy
- **Atrybut zawierający test:** Podziel wartości listy zgodnie z testem
- **Atrybut niewierający test:**
  - Nie można korzystać z informacji w funkcji testu.
  - Skorzystaj z *rid*
  - Skorzystaj z tablicy haszującej
- **Przy podziale atrybutu zawierający test:** wstaw *rid* rekordów do tablicy haszującej.
- **Tablica haszująca:** informacje o tym do którego poddrzewa rekord został przeniesiony.
- **Algorytm:**
  - Przeglądaj kolejny rekord listy
  - Dla każdego rekordu wyznacz (na podstawie tablicy haszującej) poddrzewo, do którego rekord ma być przeniesiony



# PROBLEM: ZBYT DUŻA TABLICA HASZUJĄCA

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja  
równoległa

## ALGORYTM:

- Krok 1: Podziel zbiór wartości atrybutu testującego na małe porcje tak, żeby tablica haszująca mieściła się w pamięci
- Krok 2: Dla każdej porcji
  - Podziel rekordy atrybutu testującego do właściwego podrzewa
  - Buduj tablicę haszującą
  - Przeglądaj kolejny rekord atrybutu nietestującego i przynieś go do odpowiedniego poddrzewa jeśli rekord występuje w tablicy haszującej
- Krok 3: Jeśli wszystkie rekordy zostały przydzielone do poddrzew **stop**, wpp. **idź do** krok 2



# OUTLINE

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału

SPRINT - wersja równoległa

## 1 Motywacje

## 2 ALGORYTM SPRINT

- Szukanie najlepszego podziału
- Dokonanie podziału
- SPRINT - wersja równoległa



# RÓWNOLEGŁY SPRINT

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

- Listy wartości atrybutów są równo podzielone
- **Atrybut rzeczywisty:** sortuj zbiór wartości i podziel go na równe przedziały
- **Atrybut numeryczny:** podziel według rid
- Każdy procesor ma jedną część każdej listy



# SZUKANIE NAJLEPSZEGO PODZIAŁU

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału

SPRINT - wersja równoległa

- Dla atrybutu rzeczywistego:
  - Każdy procesor ma przedział wartości atrybutu
  - Każdy procesor inicjalizuje  $C_{below}$  i  $C_{above}$  uwzględniając rozkład klas w innych procesorach
  - Każdy procesor przegląda swoją listę i wyznacza najlepszą lokalną wartość progową
  - Procesory komunikują się w celu znalezienia globalnie najlepszego cięcia
- Dla atrybutu symbolicznego:
  - Każdy procesor buduje lokalne count matrix i wysyła wynik do centralnego procesora
  - Centralny procesor oblicza globalny count matrix
  - Procesory wyznaczają najlepszy podział na podstawie globalnego count matrix



# PRZYKŁAD

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału

SPRINT - wersja równoległa

Processor 0

<i>Age</i>	<i>Class</i>	<i>rid</i>
17	High	1
20	High	5
23	High	0

<i>Car Type</i>	<i>Class</i>	<i>rid</i>
family	High	0
sports	High	1
sports	High	2

Processor 1

<i>Age</i>	<i>Class</i>	<i>rid</i>
32	Low	4
43	High	2
68	Low	3

<i>Car Type</i>	<i>Class</i>	<i>rid</i>
family	Low	3
truck	Low	4
family	high	5



# DOKONANIE PODZIAŁU

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

- Podział atrybutu zawierający test: Każdy procesor wyznacza poddrzewa, do których rekordy w lokalnej liście będą przeniesione
- Procesory wymieniają ze sobą informacje  $\langle rids, poddrzewo \rangle$
- Podział pozostałych atrybutów: Po otrzymaniu informacji ze wszystkich procesorów każdy procesor buduje tablicę haszującą i wykonuje podziały dla pozostałych atrybutów





# WADY ALGORYTMU SPRINT

Data mining

Nguyen Hung Son

Motywacje

Algorytm SPRINT

Szukanie najlepszego podziału

Dokonanie podziału  
SPRINT - wersja równoległa

- Dodatkowe struktury danych
- Nieefektywny jeśli atrybut ma dużo wartości
- Nie wykorzystuje mocnych narzędzi systemów baz danych