



Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

WYKŁAD 8: LENIWE METODY KLASYFIKACJI

Nguyen Hung Son

Wydział MIM, Uniwersytet Warszawski



OUTLINE

Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne

Szukanie minimalnych
reguł decyzyjnych

Metoda

1 LENIWE TECHNIKI KLASYFIKACJI

2 Regułowe klasyfikatory

- Reguły decyzyjne
- Szukanie minimalnych reguł decyzyjnych
- Metoda



LAZY VS. EAGER LEARNING

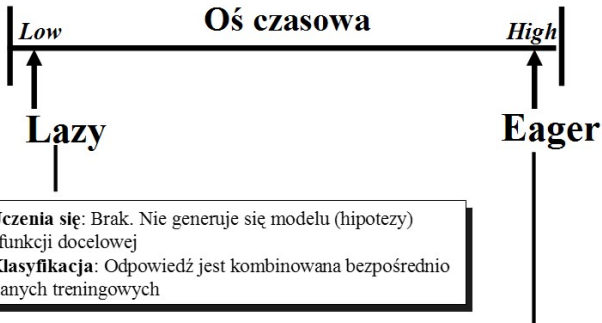
Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda





LAZY VS. EAGER LEARNING

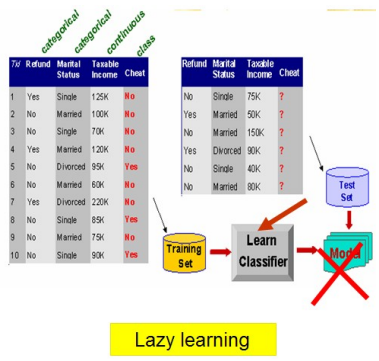
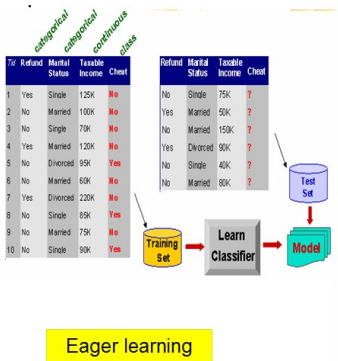
Data mining

Nguyen Hung Son

Leniwe techniki klasyfikacji

Regułowe klasyfikatory

Reguły decyzyjne
Szukanie minimalnych reguł decyzyjnych
Metoda





KIEDY STOSOWAĆ LENIWE TECHNIKI?

Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

- Eager learning: Buduje globalną hipotezę
 - Zaleta:
Prosty opis zbiorów danych
Szybki czas klasyfikacji
 - Wada:
Czas uczenia się jest obciążliwy
Jakość klasyfikacji nie wysoka
- Lazy learning: Buduje lokalną hipotezę
 - Zaleta:
Szybki czas uczenia się
Można projektować algorytmy klasyfikacji on-line
Wysoka dokładność klasyfikacji
 - Wada:
Czas klasyfikacji długi
Wymagania: funkcja odległości, strategia głosowania



Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

- k-NN
- Połączenie k-NN z drzewem decyzyjnym
- Generowanie podzbioru reguł, które mogą klasyfikować obiekt
- Generowanie poddrzewa decyzyjnego, które może klasyfikować obiekt



Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

1 Leniwe techniki klasyfikacji

2 REGUŁOWE KLASYFIKATORY

- Reguły decyzyjne
- Szukanie minimalnych reguł decyzyjnych
- Metoda



Data mining

Nguyen Hung Son

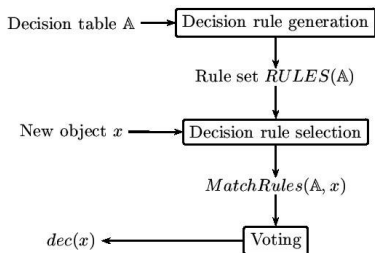
Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

Ogólny schemat:

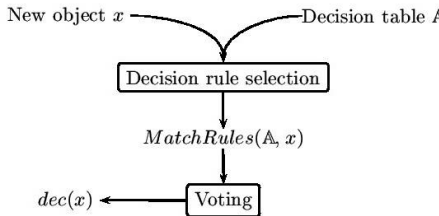
- **Uczenie się:** Generuj zbiór reguł $Rules(T)$ dla tablicy danych T
- **Selekcja reguł:** Wybierz zbiór $MatchRules(T, x)$ zawierający reguły, które pasują do nowego obiektu x .
- **klasyfikacja:** Wyznacz klasę decyzyjną dla x za pomocą głosowania na zbiorze $MatchRules(D, x)$





Ogólny schemat:

- **Generowanie reguł:**
Generuj zbiór $MatchRules(T, x)$ zawierający reguły, które pasują do nowego obiektu x .
- **klasyfikacja:** Wyznacz klasę decyzyjną dla x za pomocą głosowania na zbiorze $MatchRules(T, x)$





DESCRIPTION LANGUAGE (JĘZYK DESKRYPTORÓW)

Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

Let A be a set of attributes. The description language for A is a triple

$$\mathcal{L}(A) = (\mathbf{D}, \{\vee, \wedge, \neg\}, \mathbf{F})$$

where

- \mathbf{D} is the set of *descriptors*

$$\mathbf{D} = \{(a = v) : a \in A \text{ and } v \in Val_a\}$$

- $\{\vee, \wedge, \neg\}$ is a set of standard Boolean operators
- \mathbf{F} is a set of boolean expressions defined on \mathbf{D} called *formulas*.
- For any $B \subseteq A$ we denote by $\mathbf{D}|_B$ the set of descriptors restricted to B where
$$\mathbf{D}|_B = \{(a = v) : a \in B \text{ and } v \in Val_a\}$$
We also denote by $\mathbf{F}|_B$ the set of formulas build from $\mathbf{D}|_B$.

THE SEMANTICS

Let $\mathbb{S} = (U, A)$ be an information table describing a sample $U \subset \mathbb{X}$. The semantics of any formula $\phi \in \mathbf{F}$, denoted by $[[\phi]]_{\mathbb{S}}$, is defined by induction as follows:

$$[[a = v]]_{\mathbb{S}} = \{x \in U : a(x) = v\} \quad (1)$$

$$[[\phi_1 \vee \phi_2]]_{\mathbb{S}} = [[\phi_1]]_{\mathbb{S}} \cup [[\phi_2]]_{\mathbb{S}} \quad (2)$$

$$[[\phi_1 \wedge \phi_2]]_{\mathbb{S}} = [[\phi_1]]_{\mathbb{S}} \cap [[\phi_2]]_{\mathbb{S}} \quad (3)$$

$$[[\neg\phi]]_{\mathbb{S}} = U \setminus [[\phi]]_{\mathbb{S}} \quad (4)$$

We associate with every formula ϕ the following numeric features:

- $length(\phi)$ = the number of descriptors that occur in ϕ ;
- $support(\phi) = |[[\phi]]_{\mathbb{S}}|$ = the number of objects that match the formula;



DEFINICJA REGUŁ DECYZYJNYCH

Let $\mathbb{S} = \{U, A \cup \{dec\}\}$ be a decision table. Any implication of a form

$$\phi \Rightarrow \delta$$

where $\phi \in \mathbf{F}_A$ and $\delta \in \mathbf{F}_{dec}$, is called *the decision rule* in \mathbb{S} . The formula ϕ is called *the premise* of the decision rule \mathbf{r} and δ is called *the consequence* of \mathbf{r} . We denote the premise and the consequence of the decision rule \mathbf{r} by $prev(\mathbf{r})$ and $cons(\mathbf{r})$, respectively.



Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

GENERIC DECISION RULE

The decision rule \mathbf{r} whose the premise is a boolean monomial of descriptors, i.e.,

$$\mathbf{r} \equiv (a_{i_1} = v_1) \wedge \dots \wedge (a_{i_m} = v_m) \Rightarrow (dec = k) \quad (5)$$

is called *the generic decision rule*.

We will consider generic decision rules only. For a simplification, we will talk about decision rules keeping in mind the generic ones.



Every decision rule \mathbf{r} of the form (5) can be characterized by the following featured:

$length(\mathbf{r})$ = the number of descriptor on the assumption of \mathbf{r} (i.e. the left hand side of implication)

$[\mathbf{r}]$ = the carrier of \mathbf{r} , i.e. the set of objects from U satisfying the assumption of \mathbf{r}

$support(\mathbf{r})$ = the number of objects satisfying the assumption of \mathbf{r} : $support(\mathbf{r}) = card([\mathbf{r}])$

$$confidence(\mathbf{r}) = \frac{|[\mathbf{r}] \cap DEC_k|}{|[\mathbf{r}]|}$$

The decision rule \mathbf{r} is called *consistent* with \mathbb{A} if

$$confidence(\mathbf{r}) = 1$$



MINIMALNE REGUŁY

Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

MINIMAL CONSISTENT RULES

For a given decision table $\mathbb{S} = (U, A \cup \{dec\})$, the consistent rule:

$$\mathbf{r} = \phi \Rightarrow (dec = k)$$

is called the *minimal consistent decision rule* if any decision rule $\phi' \Rightarrow (dec = k)$ where ϕ' is a shortening of ϕ is not consistent with \mathbb{S} .



METODA OPARTA O WNIOSKOWANIE BOOLOWSKIE

Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

- Każda reguła powstaje poprzez skracanie opisu jakiegoś obiektu.
- Redukty lokalne
- Te same heurystyki dla reduktów decyzyjnych.



PRZYKŁAD TABLICY DECYZYJNEJ

Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

Hurt.	Jakość obsługi	Jakość towaru	Obok autostrady?	Centrum?	de
ID	a_1	a_2	a_3	a_4	
1	dobra	dobra	nie	nie	s
2	dobra	dobra	nie	tak	s
3	bdb	dobra	nie	nie	:
4	slaba	super	nie	nie	:
5	slaba	niska	tak	nie	:
6	slaba	niska	tak	tak	s
7	bdb	niska	tak	tak	:
8	dobra	super	nie	nie	s
9	dobra	niska	tak	nie	:
10	slaba	super	tak	nie	:
11	dobra	super	tak	tak	:
12	bdb	super	nie	tak	:
13	bdb	dobra	tak	nie	
14	slaba	super	nie	tak	



MACIERZ I FUNKCJA LOKALNYCH ODRÓŻNIALNOŚCI

Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

M	1	2	6	8
3	a_1	a_1, a_4	a_1, a_2, a_3, a_4	a_1, a_2
4	a_1, a_2	a_1, a_2, a_4	a_2, a_3, a_4	a_1
5	a_1, a_2, a_3	a_1, a_2, a_3, a_4	a_4	a_1, a_2, a_3
7	a_1, a_2, a_3, a_4	a_1, a_2, a_3	a_1	a_1, a_2, a_3, a_4
9	a_2, a_3	a_2, a_3, a_4	a_1, a_4	a_2, a_3
10	a_1, a_2, a_3	a_1, a_2, a_3, a_4	a_2, a_4	a_1, a_3
11	a_2, a_3, a_4	a_2, a_3	a_1, a_2	a_3, a_4
12	a_1, a_2, a_4	a_1, a_2	a_1, a_2, a_3	a_1, a_4

$$f_{u_3} = (\alpha_1)(\alpha_1 \vee \alpha_4)(\alpha_1 \vee \alpha_2 \vee \alpha_3 \vee \alpha_4)(\alpha_1 \vee \alpha_2) = \alpha_1$$

Reguły:

$$(a_1 = \text{bdb}) \implies \text{dec} = \text{zysk}$$



MACIERZ I FUNKCJA LOKALNYCH ODRÓŻNIALNOŚCI

Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

M	1	2	6	8
3	a_1	a_1, a_4	a_1, a_2, a_3, a_4	a_1, a_2
4	a_1, a_2	a_1, a_2, a_4	a_2, a_3, a_4	a_1
5	a_1, a_2, a_3	a_1, a_2, a_3, a_4	a_4	a_1, a_2, a_3
7	a_1, a_2, a_3, a_4	a_1, a_2, a_3	a_1	a_1, a_2, a_3, a_4
9	a_2, a_3	a_2, a_3, a_4	a_1, a_4	a_2, a_3
10	a_1, a_2, a_3	a_1, a_2, a_3, a_4	a_2, a_4	a_1, a_3
11	a_2, a_3, a_4	a_2, a_3	a_1, a_2	a_3, a_4
12	a_1, a_2, a_4	a_1, a_2	a_1, a_2, a_3	a_1, a_4

$$\begin{aligned}
 f_{u_8} &= (\alpha_1 \vee \alpha_2)(\alpha_1)(\alpha_1 \vee \alpha_2 \vee \alpha_3)(\alpha_1 \vee \alpha_2 \vee \alpha_3 \vee \alpha_4)(\alpha_2 \vee \\
 &\quad (\alpha_1 \vee \alpha_3)(\alpha_3 \vee \alpha_4)(\alpha_1 \vee \alpha_4) \\
 &= \alpha_1(\alpha_2 \vee \alpha_3)(\alpha_3 \vee \alpha_4) = \alpha_1\alpha_3 \vee \alpha_1\alpha_2\alpha_4
 \end{aligned}$$

Reguły:

- $(a_1 = \text{dobra}) \wedge (a_3 = \text{nie}) \implies \text{dec} = \text{strata}$
- $(a_1 = \text{dobra}) \wedge (a_2 = \text{super}) \wedge (a_4 = \text{nie}) \implies \text{dec} = \text{strata}$



ALGORITHM: RULE SELECTION

Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

ALGORITHM: Rule selection

Input: The object x , the maximal length λ_{\max} , the minimal support σ_{\min} , and the minimal confidence α_{\min} .

Output: The set $MatchRules(A, x)$ of decision rules from $MinRules(A, \lambda_{\max}, \sigma_{\min}, \alpha_{\min})$ matching x .

BEGIN

$C_1 := P_1; i := 1;$

WHILE ($i \leq \lambda_{\max}$) AND (C_i IS NOT EMPTY)) DO

$F_i := \emptyset; R_i := \emptyset;$

FOR $C \in C_i$ DO

$(s_1, \dots, s_d) := GetClassDistribution(C);$

$support = s_1 + \dots + s_d;$

IF $support \geq \sigma_{\min}$ THEN

IF ($\max\{s_1, \dots, s_d\} \geq \alpha_{\min} * support$) THEN

$R_i := R_i \cup \{C\};$

ELSE

$F_i := F_i \cup \{C\};$

ENDFOR

$C_{i+1} := AprGen(F_i); i := i + 1;$

ENDWHILE

RETURN $\bigcup_i R_i$

END



PRZYKŁAD

Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych

Metoda

\mathbb{A}	a_1	a_2	a_3	a_4	dec
ID	outlook	temp.	hum.	windy	play
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no
x	sunny	mild	high	TRUE	?

\Rightarrow

$\mathbb{A} _x$	d_1	d_2	d_3	d_4	dec
ID	$a_1 _x$	$a_2 _x$	$a_3 _x$	$a_4 _x$	dec
1	1	0	1	0	no
2	1	0	1	1	no
3	0	0	1	0	yes
4	0	1	1	0	yes
5	0	0	0	0	yes
6	0	0	0	1	no
7	0	0	0	1	yes
8	1	1	1	0	no
9	1	0	0	0	yes
10	0	1	0	0	yes
11	1	1	0	1	yes
12	0	1	1	1	yes
13	0	0	0	0	yes
14	0	1	1	1	no



ZBIÓR WSZTSTKICH REGUŁ

Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych
Metoda

rules	supp
outlook(overcast) \Rightarrow play(yes)	4
humidity(normal) AND windy(FALSE) \Rightarrow play(yes)	4
outlook(sunny) AND humidity(high) \Rightarrow play(no)	3
outlook(rainy) AND windy(FALSE) \Rightarrow play(yes)	3
outlook(sunny) AND temperature(hot) \Rightarrow play(no)	2
outlook(rainy) AND windy(TRUE) \Rightarrow play(no)	2
outlook(sunny) AND humidity(normal) \Rightarrow play(yes)	2
temperature(cool) AND windy(FALSE) \Rightarrow play(yes)	2
temperature(mild) AND humidity(normal) \Rightarrow play(yes)	2
temperature(hot) AND windy(TRUE) \Rightarrow play(no)	1
outlook(sunny) AND temperature(mild) AND windy(FALSE) \Rightarrow play(no)	1
outlook(sunny) AND temperature(cool) \Rightarrow play(yes)	1
outlook(sunny) AND temperature(mild) AND windy(TRUE) \Rightarrow play(yes)	1
temperature(hot) AND humidity(normal) \Rightarrow play(yes)	1

$MatchRules(\mathbb{A}, x)$ zawiera 2 reguły:

- (outlook = sunny) AND (humidity = high) \Rightarrow play = no (rule nr 3)
- (outlook = sunny) AND (temperature = mild) AND (windy = TRUE) \Rightarrow play = yes (rule nr 13)



PRZYKŁAD ALGORYTMU LENIWEGO

Data mining

Nguyen Hung Son

Leniwe techniki
klasyfikacji

Regułowe
klasyfikatory

Reguły decyzyjne
Szukanie minimalnych
reguł decyzyjnych

Metoda

$$\lambda_{max} = 3; \sigma_{min} = 1; \alpha_{min} = 1$$

$i = 1$				$i = 2$			
C_1	check	R_1	F_1	C_2	check	R_2	F_2
$\{d_1\}$	(3,2)		$\{d_1\}$	$\{d_1, d_2\}$	(1,1)	$\{d_1, d_3\}$	$\{d_1, d_2\}$
$\{d_2\}$	(4,2)		$\{d_2\}$	$\{d_1, d_3\}$	(3,0)		$\{d_2, d_3\}$
$\{d_3\}$	(4,3)		$\{d_3\}$	$\{d_1, d_4\}$	(1,1)		$\{d_1, d_4\}$
$\{d_4\}$	(3,3)		$\{d_4\}$	$\{d_2, d_3\}$	(2,2)		$\{d_2, d_3\}$
				$\{d_2, d_4\}$	(1,1)		$\{d_2, d_4\}$
				$\{d_3, d_4\}$	(2,1)		$\{d_3, d_4\}$

$i = 3$			
C_3	check	R_3	F_3
$\{d_1, d_2, d_4\}$	(0,1)	$\{d_1, d_2, d_4\}$	$\{d_2, d_3, d_4\}$
$\{d_2, d_3, d_4\}$	(1,1)		

$$MatchRules(\mathbb{A}, x) = R_2 \cup R_3:$$

(outlook = sunny) AND (humidity = high) $\Rightarrow play = no$

(outlook = sunny) AND (temperature = mild) AND (windy = TRUE)

$\Rightarrow play = yes$