

Indukcja reguł

- Kompleks k składa się z selektorów.
- $k_1 = \{< \text{słoneczna} \vee \text{deszczowa}, \text{zimna} \vee \text{ciepła}, ?, ? >\}$
 $k_2 = \{< \text{słoneczna}, \text{ciepła}, ?, ? >\}$
 $k_2 \prec k_1$
 k_2 jest bardziej szczegółowe od k_1 , k_1 jest bardziej ogólne od k_2
- $S \triangleright k$ to dokładniej $(\exists k \in S) k \triangleright x$ - zbiór wszystkich x pokrywanych przez $k \in S$
- $\{k_1 \triangleright x\} = \{1, 2, 5, 6, 9\}$
- $\{k_2 \triangleright x\} = \{1, 2\}$

Indukcja reguł - sekwencyjne pokrywanie

funkcja *sekwencyjne-pokrywanie*(T)

argumenty wejściowe:

- T - zbiór trenujący dla pojęcia c

zwraca: zbiór reguł reprezentujący hipotezę przybliżającą c

$R := 0; P := T;$

jak długo $P \neq 0$ wykonaj

$k := \text{znajdź-kompleks}(T, P);$

$d := \text{kategoria}(k, T, P);$

$R := R \cup \{k \rightarrow d\};$

$P := P - P_k;$

koniec jak długo

zwróć R

Indukcja reguł - algorytm CN2

funkcja *znajdź-kompleks-cn2*(T, P)

argumenty wejściowe:

- T - zbiór trenujący dla pojęcia c ,
- P - podzbiór zbioru T zawierający przykłady nie pokryte przez wygenerowane wcześniej reguły

zwraca: statystycznie istotny kompleks pokrywający pewną liczbę przykładów z P z dużą dokładnością;

$S := \{<? >\}; k_* := <? >;$

jak długo $S \neq \phi$ wykonaj

$S' := S \cap \mathbb{S};$

$S' := S' - S - \{< \phi >\};$

dla wszystkich kompleksów $k \in S'$ wykonaj

jeśli $\psi_k(P) > \theta \wedge \vartheta_k(P) > \vartheta_{k_*}(P)$ to $k_* := k$

koniec jeśli

koniec dla

$S := \text{Arg max}_{k \in S'}^m v_k(P)$

koniec jak długo

zwróć k_*

Algorytm CN2 - funkcja oceniająca kompleksy

Entropię zbioru P ze względu na kompleks k określa się następująco:

$$E_k(P) = \sum_{d \in C} -\frac{|P_k^d|}{|P_k|} \log \frac{|P_k^d|}{|P_k|}$$

Entropia ma tę cechę, że największą wartość przyjmuje dla zrównoważonych rozkładów częstości kategorii. Funkcja oceniająca kompleksy musi być zane-gowaną entropią:

$$\vartheta_k(P) = -E_k(P)$$

Algorytm CN2 - statystyka χ

Niech f_i oznacza *zaobserwowaną częstość* (liczbę wystąpień) i -tej wartości atrybutu y_i dla $i = 1, 2, 3, \dots, v_1$ i odpowiednio f_j dla y_j dla $j = 1, 2, 3, \dots, v_2$, f_{ij} liczbę (częstość) jednoczesnych wystąpień i -tej i j -tej wartości atrybutów y_i i y_j , a e_{ij} to wartość oczekiwana jednoczesnego wystąpienia przy założeniu niezależności y_1 i y_2 i $(v_1 - 1)(v_2 - 1)$ stopniach swobody.

$$\chi^2 = \sum_{i=1}^{v_1} \sum_{j=1}^{v_2} \frac{(f_{ij} - e_{ij})^2}{e_{ij}},$$

gdzie $e_{ij} = \frac{f_i^1 f_j^2}{n}$

Im większa wartość statystyki tym bardziej atrybuty są zależne od siebie.

Algorytm CN2 - statystyka χ

$$\chi_k^2(P) = \sum_{d \in C} \frac{(|P_k^d| - e_k^d(P))^2}{e_k^d(P)},$$

$$\text{gdzie } e_k^d(P) = |P_k| \frac{|P^d|}{|P|}$$

x	<i>aura</i>	<i>temperatura</i>	<i>wilgotność</i>	<i>wiatr</i>	$c(x)$
1	<i>słoneczna</i>	<i>ciepła</i>	<i>duża</i>	<i>słaby</i>	0
2	<i>słoneczna</i>	<i>ciepła</i>	<i>duża</i>	<i>silny</i>	0
3	<i>pochmurna</i>	<i>ciepła</i>	<i>duża</i>	<i>słaby</i>	1
4	<i>deszczowa</i>	<i>umiarkowana</i>	<i>duża</i>	<i>słaby</i>	1
5	<i>deszczowa</i>	<i>zimna</i>	<i>normalna</i>	<i>słaby</i>	1
6	<i>deszczowa</i>	<i>zimna</i>	<i>normalna</i>	<i>silny</i>	0
7	<i>pochmurna</i>	<i>zimna</i>	<i>normalna</i>	<i>silny</i>	1
8	<i>słoneczna</i>	<i>umiarkowana</i>	<i>duża</i>	<i>słaby</i>	0
9	<i>słoneczna</i>	<i>zimna</i>	<i>normalna</i>	<i>słaby</i>	1
10	<i>deszczowa</i>	<i>umiarkowana</i>	<i>normalna</i>	<i>słaby</i>	1
11	<i>słoneczna</i>	<i>umiarkowana</i>	<i>normalna</i>	<i>silny</i>	1
12	<i>pochmurna</i>	<i>umiarkowana</i>	<i>duża</i>	<i>silny</i>	1
13	<i>pochmurna</i>	<i>ciepła</i>	<i>normalna</i>	<i>słaby</i>	1
14	<i>deszczowa</i>	<i>umiarkowana</i>	<i>duża</i>	<i>silny</i>	0

Zbiór \mathbb{S} kompleksów atomowych

$\mathbb{S} = \{ \langle \text{deszczowa}, ?, ?, ? \rangle,$
 $\langle \text{deszczowa} \vee \text{słoneczna}, ?, ?, ? \rangle,$
 $\langle \text{deszczowa} \vee \text{pochmurna}, ?, ?, ? \rangle,$
 $\langle \text{pochmurna}, ?, ?, ? \rangle,$
 $\langle \text{pochmurna} \vee \text{słoneczna}, ?, ?, ? \rangle,$
 $\langle \text{słoneczna}, ?, ?, ? \rangle,$
 $\langle ?, \text{ciepła}, ?, ? \rangle,$
 $\langle ?, \text{ciepła} \vee \text{zimna}, ?, ? \rangle,$
 $\langle ?, \text{ciepła} \vee \text{umiarkowana}, ?, ? \rangle,$
 $\langle ?, \text{umiarkowana}, ?, ? \rangle,$
 $\langle ?, \text{umiarkowana} \vee \text{zimna}, ?, ? \rangle,$
 $\langle ?, \text{zimna}, ?, ? \rangle,$
 $\langle ?, ?, \text{duża}, ? \rangle, \langle ?, ?, \text{normalna}, ? \rangle, \langle ?, ?, ?, \text{silny} \rangle, \langle ?, ?, ?, \text{słaby} \rangle \}$

1. Początkowo $R = \phi, P = T = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14\}, \mathbb{S}$
2. Następuje wywołanie *znajdź-kompleks*(T, P).
 - $S = \{<? >\} \neq \phi, k_* = <? >$ i $\vartheta_{k_*}(P) = -E_{k_*}(P) = -0.940$,
 - $S' = \mathbb{S} = S \cap \mathbb{S}$,
 - $k = < \text{pochmurna}, ?, ?, ? >$ ma największą wartość $\vartheta_k = 0$ w zbiorze \mathbb{S} ;
 $S = \{k\}, k_* = k$,
3. $R = \{< \text{pochmurna}, ?, ?, ? > \rightarrow 1\}, P = \{1, 2, 4, 5, 6, 8, 9, 10, 11, 14\}$,
4. $P \neq \phi \Rightarrow \text{znajdź-kompleks}(T, P)$,
 - $S = \{<? >\} \neq \phi, k_* = <? >$ i $\vartheta_{k_*}(P) = -1$,
 - $S' = \mathbb{S} = S \cap \mathbb{S}$,
 - $k = <?, \text{ciepła}, ?, ? >$ ma największą wartość $\vartheta_k = 0$ w zbiorze \mathbb{S} ;
 $S = \{k\} \neq \phi, k_* = k$,
5. $R = \{< \text{pochmurna}, ?, ?, ? > \rightarrow 1, <?, \text{ciepła}, ?, ? > \rightarrow 0\}$,
 $P = \{4, 5, 6, 8, 9, 10, 11, 14\}$,

1. $P \neq \phi \Rightarrow \text{znajdź-kompleks}(T, P)$,
 - $S' = \mathbb{S} = S \cap \mathbb{S}$,
 - $k = \langle ?, ?, \text{normalna}, ? \rangle$ zostaje wybrane z najwyższą wartością $\vartheta_k = -0,721$ w zbiorze \mathbb{S} ; $S = \{k\} \neq \phi$, $k_* = k$,
 - k_* nie ma wartości 0 (pętla jak długo się nie kończy),
 - w następnym cyklu dla $S' = S \cap \mathbb{S}$ największą wartość $\vartheta_k = 0$ ma kompleks $k = \langle ?, ?, \text{normalna}, \text{słaby} \rangle$, $k_* = k$
2. $R = \{ \langle \text{pochmurna}, ?, ?, ? \rangle \rightarrow 1, \langle ?, \text{ciepła}, ?, ? \rangle \rightarrow 0, \langle ?, ?, \text{normalna}, \text{słaby} \rangle \rightarrow 1 \}$, $P = \{4, 6, 8, 11, 14\}$,
3. po kilku dalszych wywołaniach funkcji $\text{znajdź-kompleks}(T, P)$ otrzymujemy
 $R = \{ \langle \text{pochmurna}, ?, ?, ? \rangle \rightarrow 1, \langle ?, \text{ciepła}, ?, ? \rangle \rightarrow 0, \langle ?, ?, \text{normalna}, \text{słaby} \rangle \rightarrow 1, \langle ?, \text{zimna}, ?, ? \rangle \rightarrow 0, \langle ?, ?, \text{normalna}, ? \rangle \rightarrow 1, \langle ?, ?, ?, \text{silny} \rangle \rightarrow 0, \langle \text{słoneczna}, ?, ?, ? \rangle \rightarrow 0 \}$, $P = \{4\}$

1. $P \neq \phi \Rightarrow \text{znajdź-kompleks}(T, P)$,
 - $S = \{ \langle ? \rangle \} \neq \phi, k_* = \langle ? \rangle$ i $\vartheta_{k_*}(P) = -E_{k_*}(P) = 0$,
2. Ostatecznie

$$R = \{ \langle \text{pochmurna}, ?, ?, ? \rangle \rightarrow 1, \\ \langle ?, \text{ciepła}, ?, ? \rangle \rightarrow 0, \\ \langle ?, ?, \text{normalna}, \text{słaby} \rangle \rightarrow 1, \\ \langle ?, \text{zimna}, ?, ? \rangle \rightarrow 0, \\ \langle ?, ?, \text{normalna}, ? \rangle \rightarrow 1, \\ \langle ?, ?, ?, \text{silny} \rangle \rightarrow 0, \\ \langle \text{słoneczna}, ?, ?, ? \rangle \rightarrow 0, \\ \langle ? \rangle \rightarrow 1 \}$$