

Dokumentacja wstępna

UMA

Andrii Gamalii,
Wiktor Topolski

Treść zadania

Algorytm analogiczny do lasu losowego z użyciem zbiorów reguł zamiast drzew decyzyjnych.

Interpretacja zadania

Las losowy z użyciem zbiorów reguł **uporządkowanych** zamiast drzew decyzyjnych.
Indukcja reguł za pomocą specjalizacji AQ.

Opis algorytmów

Zmodyfikowany las losowy

B zbiorów reguł (oryginalnie drzew decyzyjnych).

Uczenie

Hiperparametry

- B - maksymalna liczba zbiorów reguł
- M - wielkość podzbioru trenującego dla każdej reguły

Opis algorytmu

Nauczenie B zbiorów reguł na losowo wybranym podzbiorze zbioru trenującego o wielkości M (nie usuwając tych przykładów ze zbioru trenującego) (nie może mieć sprzecznych przykładów), z losowo (bez zwracania) wybranymi atrybutami o liczbie równej podłódze pierwiastka maksymalnej liczby atrybutów. Sposób nauczania (indukcja reguł) - algorytm AQ - opisany niżej.

Predykcja

Parametry

- B nauczonych zbiorów reguł
- V - wektor atrybutów

Opis algorytmu

Wybieramy tą klasę, którą wybrała większość zbiorów reguł.

Specjalizacja AQ

Algorytm tworzący zbiór reguł dla danego zbioru trenującego.

Algorytm zakłada dyskretne i skończone zbiory wartości atrybutów.

- **Kompleks** - $\langle s_1, s_2, \dots, s_n \rangle$ - wektor n (koniunkcja) selektorów dla poszczególnych n atrybutów, spełniony gdy wszystkie selektory są spełnione, występują różne warunki np. $\langle 2, 3 \vee 2, ? \rangle$
? - selektor uniwersalny
- **Reguła** - kompleks i klasa którą mają przykłady spełniane przez kompleks.
- **Zbiór reguł** (uporządkowany) - każda reguła stosowana tylko dla przykładów niepokrytych przez wcześniejsze reguły, ostatnia reguła określa klasę domyślną, brak możliwości konfliktu reguł o różnych klasach.

Uczenie

Hiperparametry

- m - Maksymalna ilość kompleksów, które pozostają po specjalizacji (najlepiej ocenione kompleksy)
- T - Maksymalna ilość reguł w jednym zbiorze reguł

Modyfikowalne składowe algorytmu

- Sposób wyboru ziarna
- Sposób oceny kompleksu

Opis algorytmu

Algorytm powtarzamy aż wszystkie przykłady zostaną pokryte, lub gdy osiągniemy maksymalną liczbę reguł w zbiorze (t).

Ziarno x_s - przykład podstawowy, pierwszy wybierany, jest zawsze pokryty przez nową regułę.

$R^{(0)}_G$ - zbiór wszystkich pokrywanych przykładów niewłaściwej klasy przez kompleksy ze zbioru G w danej iteracji.

R - niepokryte przykłady

$G = \langle ? \rangle$;

$x_s \in R$ - ziarno

$R^{(1)} = R_{c=c(x_s)}$ - klasa ziarna

$R^{(0)} = R_{c \neq c(x_s)}$

jak długo $R^{(0)}_G \neq \emptyset$ (kiedy są rzeczy w tym r_0):

- wybierz $x_n \in R^{(0)}_G$ - ziarno negatywne
- dla wszystkich $k \in \{k' \in G \mid k' \triangleright x_n\}$
 - $G = G - \{k\} \cup \text{specjalizacja}(k, x_n, x_s)$
- $G = G - \{k \in G \mid (\exists k' \in G) k' \succ k\}$
- $G = \text{Argmax}_{k \in G \cup \nu_{R^{(1)}, R^{(0)}}}(k)$ (m najlepszych kompleksów)

zwrot $\text{arg max}_{k \in G \cup \nu_{R^{(1)}, R^{(0)}}}(k)$ (najlepszy kompleks z G)

gdzie

$\nu_{R^{(1)}, R^{(0)}}(k)$ to ocena jakości kompleksu k

- zazwyczaj poprzez liczbę wcześniej niepokrywanych przykładów, pokrywanych przez dany kompleks

Uporządkowanie zbioru reguł:

- eliminacja pokrywania przykładów klas innych niż $c(x_s)$ wyłącznie w zbiorze R (przykładów niepokrytych przez wcześniejsze reguły)

$\text{specjalizacja}(k, x_n, x_s)$

chcemy uzyskać zbiór maksymalnie ogólnych kompleksów k' spełniających warunki:

- $k' \prec k$
- $k' \not\triangleright x_n$
- $k' \triangleright x_s$

Przykład działania

Dany zbiór trenujący R:

x	a ₁	a ₂	a ₃	c
1	1	1	3	1
2	2	2	2	1
3	3	2	1	0
4	1	1	2	0

Gdzie a₁: X → {1, 2, 3}, a₂: X → {1, 2}, a₃: X → {1, 2, 3}, c: X → {0, 1}.

Maksymalny rozmiar **m** zbioru kompleksów **G** = 2.

Proces indukcji pierwszej reguły:

1. Jako ziarno **x_s** wybieramy pierwszy przykład, czyli **x₁**.

Ziarno jest klasy 1, więc **R⁽¹⁾** = {x₁, x₂}, **R⁽⁰⁾** = {x₃, x₄}.

Licznik reguł t = 0.

Zbiór kompleksów **G** inicjalizujemy zbiorem zawierającym jeden kompleks uniwersalny:

G = {< ? >}.

Kroki 2, 3, 4, 5, 6 są dokonywane w pętli dopóki **R⁽⁰⁾_G** nie jest pusty.

2. **R⁽⁰⁾_G** = {x₃, x₄} - zbiór pokrywanych przykładów niewłaściwej klasy.
3. Jako ziarno negatywne **x_n** wybieramy x₃.
4. Kompleks **k** = < ? > jest poddawany operacji **specjalizacja(k, x_n, x_s)**, w wyniku której < ? > przechodzi na {< 1 ∨ 2, ?, ? >, < ?, 1, ? >, < ?, ?, 2 ∨ 3 >}.
5. Ze zbioru **G** usuwany jest kompleks < ? > i wstawiane są kompleksy otrzymane po jego specjalizacji. **G** = {< 1 ∨ 2, ?, ? >, < ?, 1, ? >, < ?, ?, 2 ∨ 3 >}.
6. Wielkość zbioru **G** wynosi 3, podczas gdy **m** = 2, więc musimy zostawić w nim tylko 2 najlepsze kompleksy, oceniane przez np. **pokrycie**. W tym przypadku według pokrycia najlepsze są kompleksy < 1 ∨ 2, ?, ? >, < ?, ?, 2 ∨ 3 >, więc kompleks < ?, 1, ? > jest odrzucany.

2. **G** = {< 1 ∨ 2, ?, ? >, < ?, ?, 2 ∨ 3 >}, **R⁽⁰⁾_G** = {x₄}.

3. Jako ziarno negatywne **x_n** wybieramy x₄.

4. specjalizacja(< 1 ∨ 2, ?, ? >, x₄, x₁):

< 1 ∨ 2, ?, ? > → < 1 ∨ 2, ?, 2 ∨ 3 >

specjalizacja(< ?, 1, ? >, x₄, x₁):

< ?, ?, 2 ∨ 3 > → < ?, ?, 3 >

5. **G** = {< 1 ∨ 2, ?, 2 ∨ 3 >, < ?, ?, 3 >}

6. Rozmiar zbioru **G** nie wykracza poza limit
7. $R^{(0)}_G = \{\emptyset\}$, co oznacza, że nie zostało pokrywanych przykładów niewłaściwej klasy i możemy wyjść z pętli
8. Wybieramy najlepszy kompleks ze zbioru **G**, np. według miary pokrycia najlepszy jest kompleks $\langle 1 \vee 2, ?, 2 \vee 3 \rangle$.
9. Zwracamy regułę $\langle 1 \vee 2, ?, 2 \vee 3 \rangle \rightarrow 1$.
10. $R = R \setminus \{x_i\}$, $t += 1$.
11. Powyższe kroki są powtarzane dopóki $R \neq \{\emptyset\}$ oraz $t \leq T$.

Predykcja

Pierwsza reguła spełniająca dany przykład wyznacza klasę, ostatnia reguła wyznacza klasę domyślną.

Zbiór danych

<http://archive.ics.uci.edu/dataset/73/mushroom>

Klasyfikacja czy grzyb jest trujący na podstawie 20 atrybutów dyskretnych

8124 przykłady - 4208 jadalnych grzybów, 3916 niejadalnych

Zbiór trenujący wybieramy losowo

Eksperymenty

Zmienne algorytmów

- Sposoby oceny kompleksów, np. pokrycie, dokładność, dominacja klasy z uwzględnieniem pokrycia.
- Sposoby wyboru ziarna np. losowy, pierwszy przykład, losowy z klasy dominującej.

Hiperparametry

Użyjemy przeszukiwania **losowego** wartości hiperparametrów.

- Rozmiar zbioru trenującego (dla całego algorytmu), np. 80%, 90%, 95%
- Rozmiar zbiorów trenujących dla poszczególnych zbiorów reguł, np. 500 przykładów
- Maksymalna ilość zbiorów reguł, np. 100, 500, 1400
- Maksymalna ilość reguł w zbiorze reguł, np. 2, 10, 250

Miary jakości

- Confusion matrix
- Dokładność
- Precyzja
- F1-score