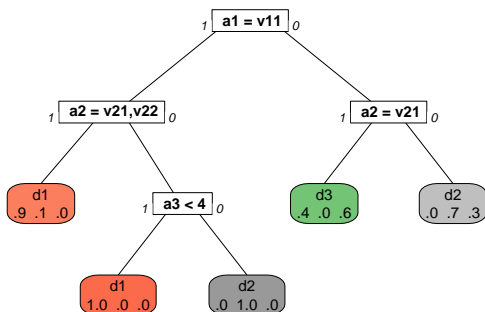


Uczenie maszynowe:
wykład 8

Paweł Cichosz

- 1 Reprezentacja modeli za pomocą drzew decyzyjnych
- 2 Budowa drzewa

Struktura drzewa



Struktura drzewa

Węzły: podziały (testy) na podstawie warunków dotyczących wartości atrybutów $t : X \rightarrow R_t$ (pełnią analogiczną rolę jak selektory w regułach).

Gałęzie: dla każdego wyniku $r \in R_t$ podziału t w węźle prowadzą z tego węzła do jego węzłów potomnych.

Liście: klasy i prawdopodobieństwa klas.

Proces predykcji: przykład x propagowany od korzenia drzewa do liścia ścieżką wyznaczaną przez wyniki podziałów w kolejnych odwiedzonych węzłach $t_1(x), t_2(x), \dots$, klasa z osiągniętego liścia jest wartością $h(x)$.

Podziały dla atrybutów dyskretnych

Wielowartościowe na podstawie wartości atrybutu:

$$t(x) = a(x)$$



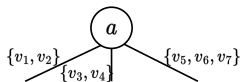
Binarne na podstawie równości:

$$t(x) = \begin{cases} 1 & \text{jeżeli } a(x) = v \\ 0 & \text{w przeciwnym przypadku} \end{cases}$$



Wielowartościowe na podstawie podzbiorów wartości atrybutu:

$$t(x) = \begin{cases} 1 & \text{jeżeli } a(x) \in V_1 \\ 2 & \text{jeżeli } a(x) \in V_2 \\ \dots & \\ k & \text{jeżeli } a(x) \in V_k \end{cases}$$



Binarne na podstawie przynależności do zbioru:

$$t(x) = \begin{cases} 1 & \text{jeżeli } a(x) \in V \\ 0 & \text{w przeciwnym przypadku} \end{cases}$$



Podziały dla atrybutów ciągłych

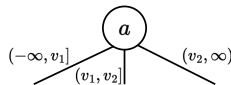
Binarne na podstawie nierówności:

$$t(x) = \begin{cases} 1 & \text{jeżeli } a(x) \leq v \\ 0 & \text{w przeciwnym przypadku} \end{cases}$$

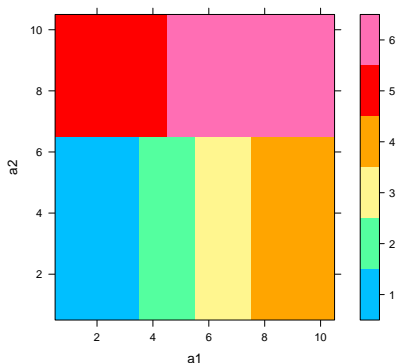
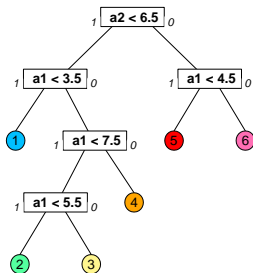


Wielowartościowe na podstawie przedziałów wartości atrybutu:

$$t(x) = \begin{cases} 1 & \text{jeżeli } a(x) \in (-\infty, v_1] \\ 2 & \text{jeżeli } a(x) \in (v_1, v_2] \\ \dots & \\ k & \text{jeżeli } a(x) \in (v_{k-1}, \infty) \end{cases}$$



Dekompozycja dziedziny



Z każdym węzłem lub liściem n jest związany odpowiedni podzbiór X_n dziedziny X i odpowiedni podzbiór S_n dowolnego zbioru danych S , zawierający przykłady, które docierają do tego węzła lub liścia, jeśli są propagowane z korzenia drzewa ścieżkami wyznaczanymi przez wyniki kolejnych podziałów.

Wymiar VC

Teoretyczny, bez atrybutów ciągłych: liczba możliwych kombinacji wartości atrybutów.

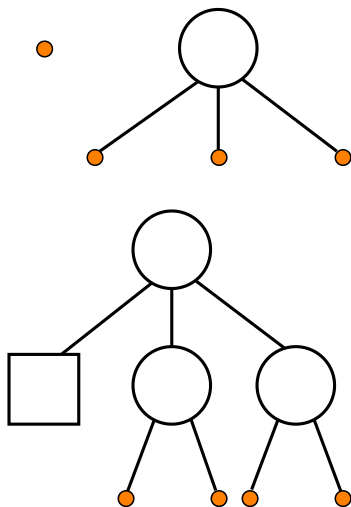
Teoretyczny, z atrybutami ciągłymi: ∞ .

Efektywny: zredukowany przez zastosowanie kryteriów stopu i przycinania (jak zobaczymy dalej).

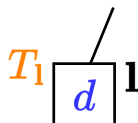
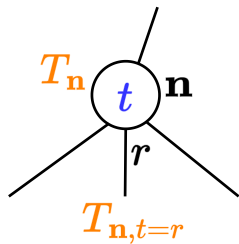
- 1 Reprezentacja modeli za pomocą drzew decyzyjnych
- 2 Budowa drzewa

Schemat algorytmu

- Sekwencja decyzji:
 - kryterium stopu: węzeł czy liść?
 - wybór klasy dla liścia: jaka predykcja najlepsza w liściu?
 - wybór podziału dla węzła: jaki podział najlepszy w węźle?
- Ten sam schemat powtarzany dla korzenia drzewa, jego węzłów potomnych, ich węzłów potomnych itd.
- Kolejność zazwyczaj nieistotna (zwykle: w głąb albo w szerz).



Drzewa decyzyjne: notacja



Kryterium stopu

Podejmowana decyzja: czy węzeł n , któremu odpowiada podzbiór przykładów trenujących T_n , powinien być liściem?

Jednolita klasa: T_n zawiera przykłady dokładnie jednej klasy (wybieranej jako klasa dla liścia).

Brak przykładów: $T_n = \emptyset$ (klasa dla liścia z węzła macierzystego).

Brak możliwości podziału: każdy możliwy podział osiąga dla T_n tylko jeden wynik (nie dzieli).

Warianty rozluźnione: dominacja jednej klasy, mała liczba przykładów, najlepszy podział zbyt słaby.

Wpływ na właściwości drzewa: kryteria stopu (w nierozluźnionej postaci) gwarantują minimalizację błędu na zbiorze trenującym.

Uzupełniające kryterium: maksymalna liczba poziomów.

Prawdopodobieństwa klas: na podstawie rozkładu przykładów trenujących w liściu.

Kryterium wyboru podziału

Motywacja: preferencja dla prostych drzew ograniczająca ryzyko nadmiernego dopasowania.

Realizacja: ocena jakości podziałów na podstawie nieczystości rozkładu klas po podziale, wybór najlepszego podziału w każdym węźle (ze skończonego zbioru kandydatów).

Ocena jakości podziału:

- węzeł \mathbf{n} ,
- zbiór przykładów trenujących $T_{\mathbf{n}}$,
- rozważany podział $t : X \rightarrow R_t$,
- dla każdego wyniku podziału $r \in R_t$ odpowiedni podzbiór przykładów trenujących $T_{\mathbf{n},t=r} = \{x \in T_{\mathbf{n}} \mid t(x) = r\}$.

Kryterium wyboru podziału

Entropia warunkowa:

$$E_{T_{\mathbf{n}}}(c|t) = \sum_{r \in R_t} \frac{|T_{\mathbf{n},t=r}|}{|T_{\mathbf{n}}|} E_{T_{\mathbf{n},t=r}}(c)$$

gdzie:

$$E_{T_{\mathbf{n},t=r}}(c) = \sum_{d \in C} -P_{T_{\mathbf{n},t=r}}(c = d) \log P_{T_{\mathbf{n},t=r}}(c = d)$$

$$P_{T_{\mathbf{n},t=r}}(c = d) = \frac{|T_{\mathbf{n},t=r,c=d}|}{|T_{\mathbf{n},t=r}|}$$

Redukcja nieczystości w wyniku podziału (przyrost informacji):

$$\Delta E_{T_{\mathbf{n}}}(c|t) = E_{T_{\mathbf{n}}}(c) - E_{T_{\mathbf{n}}}(c|t)$$

Wybór podziału: minimalizacja nieczystości po podziale (równoważnie: maksymalizacja redukcji nieczystości w wyniku podziału).

Kryterium wyboru podziału

Warunkowy indeks Giniego:

$$G_{T_n}(c|t) = \sum_{r \in R_t} \frac{|T_{n,t=r}|}{|T_n|} G_{T_{n,t=r}}(c)$$

gdzie:

$$G_{T_{n,t=r}}(c) = 1 - \sum_{d \in C} P_{T_{n,t=r}}(c = d)^2$$

Redukcja nieczystości w wyniku podziału:

$$\Delta G_{T_n}(c|t) = G_{T_n}(c) - G_{T_n}(c|t)$$

Wybór podziału: minimalizacja nieczystości po podziale (równoważnie: maksymalizacja redukcji nieczystości w wyniku podziału).

Przykład: pogoda

x	<i>outlook</i>	<i>temperature</i>	<i>humidity</i>	<i>wind</i>	<i>play</i>
1	<i>sunny</i>	<i>hot</i>	<i>high</i>	<i>normal</i>	<i>no</i>
2	<i>sunny</i>	<i>hot</i>	<i>high</i>	<i>high</i>	<i>no</i>
3	<i>overcast</i>	<i>hot</i>	<i>high</i>	<i>normal</i>	<i>yes</i>
4	<i>rainy</i>	<i>mild</i>	<i>high</i>	<i>normal</i>	<i>yes</i>
5	<i>rainy</i>	<i>cold</i>	<i>normal</i>	<i>normal</i>	<i>yes</i>
6	<i>rainy</i>	<i>cold</i>	<i>normal</i>	<i>high</i>	<i>no</i>
7	<i>overcast</i>	<i>cold</i>	<i>normal</i>	<i>high</i>	<i>yes</i>
8	<i>sunny</i>	<i>mild</i>	<i>high</i>	<i>normal</i>	<i>no</i>
9	<i>sunny</i>	<i>cold</i>	<i>normal</i>	<i>normal</i>	<i>yes</i>
10	<i>rainy</i>	<i>mild</i>	<i>normal</i>	<i>normal</i>	<i>yes</i>
11	<i>sunny</i>	<i>mild</i>	<i>normal</i>	<i>high</i>	<i>yes</i>
12	<i>overcast</i>	<i>mild</i>	<i>high</i>	<i>high</i>	<i>yes</i>
13	<i>overcast</i>	<i>hot</i>	<i>normal</i>	<i>normal</i>	<i>yes</i>
14	<i>rainy</i>	<i>mild</i>	<i>high</i>	<i>high</i>	<i>no</i>

Przykład: pogoda

x	<i>outlook</i>	<i>temperature</i>	<i>humidity</i>	<i>wind</i>	<i>play</i>
1	<i>sunny</i>	<i>hot</i>	<i>high</i>	<i>normal</i>	<i>no</i>
2	<i>sunny</i>	<i>hot</i>	<i>high</i>	<i>high</i>	<i>no</i>
3	<i>overcast</i>	<i>hot</i>	<i>high</i>	<i>normal</i>	<i>yes</i>
4	<i>rainy</i>	<i>mild</i>	<i>high</i>	<i>normal</i>	<i>yes</i>
5	<i>rainy</i>	<i>cold</i>	<i>normal</i>	<i>normal</i>	<i>yes</i>
6	<i>rainy</i>	<i>cold</i>	<i>normal</i>	<i>high</i>	<i>no</i>
7	<i>overcast</i>	<i>cold</i>	<i>normal</i>	<i>high</i>	<i>yes</i>
8	<i>sunny</i>	<i>mild</i>	<i>high</i>	<i>normal</i>	<i>no</i>
9	<i>sunny</i>	<i>cold</i>	<i>normal</i>	<i>normal</i>	<i>yes</i>
10	<i>rainy</i>	<i>mild</i>	<i>normal</i>	<i>normal</i>	<i>yes</i>
11	<i>sunny</i>	<i>mild</i>	<i>normal</i>	<i>high</i>	<i>yes</i>
12	<i>overcast</i>	<i>mild</i>	<i>high</i>	<i>high</i>	<i>yes</i>
13	<i>overcast</i>	<i>hot</i>	<i>normal</i>	<i>normal</i>	<i>yes</i>
14	<i>rainy</i>	<i>mild</i>	<i>high</i>	<i>high</i>	<i>no</i>

Wybór podziału dla korzenia drzewa, w którym
 $T_n = T$.

- Dla podziału *outlook* (stosując logarytmy dwójkowe):

$$E_{T_{outlook=sunny}}(c) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \approx 0.971$$

$$E_{T_{outlook=overcast}}(c) = -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} = 0$$

$$E_{T_{outlook=rainy}}(c) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \approx 0.971$$

- Entropia warunkowa:

$$E_T(c|outlook) = \frac{5}{14} \cdot 0.971 + \frac{4}{14} \cdot 0 + \frac{5}{14} \cdot 0.971 \approx 0.694$$

- Analogicznie dla pozostałych podziałów (kontynuacja pozostaje jako ćwiczenie).