

# Uczenie maszynowe: *wykład 9*

Paweł Cichosz

- 1 Przycinanie drzew decyzyjnych
- 2 Naiwny klasyfikator bayesowski

# Przycinanie drzewa

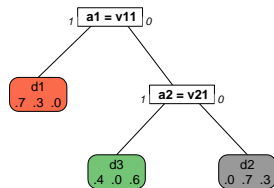
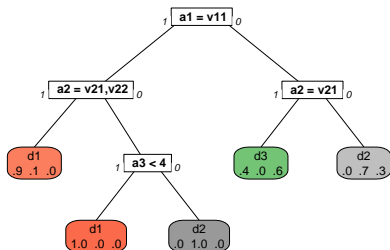
**Cel:** Zapobieganie nadmiernemu dopasowaniu (bardziej wymagające, ale potencjalnie skuteczniejsze niż kryterium stopu).

**Operator:** zastąpienie poddrzewa liściem.

**Kryterium:** oczekiwana redukcja błędu rzeczywistego (wiele różnych szczegółowych wariantów).

**Strategia:** najczęściej wstępująca (czasem zstępująca, najpierw najlepszy).

## Przycinanie drzewa



## REP (Reduced Error Pruning)

- Kryterium przycinania oparte na estymacji błędu rzeczywistego z wykorzystaniem osobnego podzbioru przykładów  $R$  nieużywanego w trakcie budowy drzewa.
- Węzeł  $\mathbf{n}$  zastępowany liściem  $\mathbf{l}$  jeśli:

$$e_R(\mathbf{l}) \leq e_R(\mathbf{n})$$

przy czym błąd liścia  $e_R(\mathbf{l})$  i błąd węzła  $e_R(\mathbf{n})$  wyznacza się na podstawie tych przykładów z  $R$ , które „docierają” do miejsca jego położenia w drzewie.

- Dobre jeśli możemy poświęcić osobną pulę przykładów.
- Może czasem lepiej przycinać trochę gorzej, ale budować lepiej (ale inne metody poza zakresem wykładu)?

## Konwersja drzewa do zbioru reguł

- Reguła dla każdej ścieżki od korzenia do liścia:  
część warunkowa: koniunkcja warunków odpowiadających wynikom podziałów w węzłach,  
część decyzyjna: klasa z liścia.
- Reprezentacja uważana (czasem) za bardziej czytelną dla człowieka.
- Umożliwia bardziej elastyczne przycinanie (usuwanie bądź pozostawianie warunków niezależnie dla każdej reguły).

## Właściwości drzew decyzyjnych

- Zwykle dobra jakość predykcji, chociaż często inne algorytmy dają nieco lepsze modele.
- Reprezentacja modeli czytelna dla człowieka.
- Podatne na nadmierne dopasowanie – konieczne staranne dobranie kryteriów stopu lub zastosowanie przycinania.
- Po niewielkiej modyfikacji mogą służyć jako modele regresji (poza zakresem wykładu).
- Przydatne jako komponenty modeli zespołowych, takich jak las losowy (omawiany na jednym z kolejnych wykładów) i *gradient boosting* (poza zakresem wykładu).

- 1 Przycinanie drzew decyzyjnych
- 2 Naiwny klasyfikator bayesowski



# Naiwny klasyfikator bayesowski

Wzór Bayesa:

$$P(A|B) = \frac{P(A)P(B|A)}{P(B)}$$

Zastosowanie do predykcji prawdopodobieństw klas:

$$\pi_d(x) = P(d|x) = P(c = d \mid a_1 = a_1(x), a_2 = a_2(x), \dots, a_n = a_n(x))$$

Prawdopodobieństwo klasy na podstawie wartości atrybutów:

$$\begin{aligned} P(c = d \mid a_1 = v_1, \dots, a_n = v_n) \\ = \frac{P(c = d)P(a_1 = v_1, \dots, a_n = v_n \mid c = d)}{P(a_1 = v_1, \dots, a_n = v_n)} \end{aligned}$$

Założenie o niezależności:

$$P(a_1 = v_1, \dots, a_n = v_n \mid c = d) = \prod_{i=1}^n P(a_i = v_i \mid c = d)$$

# Estymacja prawdopodobieństw

Prawdopodobieństwo *a priori* klasy:

$$P(c = d) = P_T(c = d) = \frac{|T_{c=d}|}{|T|}$$

Prawdopodobieństwo warunkowe wartości atrybutu:

$$P(a_i = v_i | c = d) = P_{T_{c=d}}(a_i = v_i) = \frac{|T_{c=d, a_i=v_i}|}{|T_{c=d}|}$$

Mianownik: stała normalizująca:

$$P(a_1 = v_1, \dots, a_n = v_n) = \sum_{d \in C} P(c = d) P(a_1 = v_1, \dots, a_n = v_n | c = d)$$

## Prawdopodobieństwa zerowe i prawie zerowe

Stała wartość dodatnia:  $P(a_i = v_i | c = d) = \epsilon$  jeśli  $T_{c=d, a_i=v_i} = \emptyset$ .

Wygładzanie Laplace'a:

$$P(a_i = v_i | c = d) = \frac{|T_{c=d, a_i=v_i}| + l}{|T_{c=d}| + l|A_i|}$$

$$(l \geq 0)$$

**Logarytm prawdopodobieństwa:** mniejsze ryzyko błędów numerycznych przy mnożeniu małych prawdopodobieństw:

$$\begin{aligned} \log \left( P(c = d) \prod_{i=1}^n P(a_i = v_i | c = d) \right) \\ = \log P(c = d) + \sum_{i=1}^n \log P(a_i = v_i | c = d) \end{aligned}$$

## 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

<i>x</i>	<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>

Klasyfikacja przykładu:

 15 | *sunny* *mild* *normal* *normal* | *no*

$$P(\text{no}) = \frac{5}{14}$$

$$P(\text{yes}) = \frac{9}{14}$$

$$P(\text{outlook} = \text{sunny}|\text{no}) = \frac{3}{5}$$

$$P(\text{outlook} = \text{sunny}|\text{yes}) = \frac{2}{9}$$

$$P(\text{temperature} = \text{mild}|\text{no}) = \frac{2}{5}$$

$$P(\text{temperature} = \text{mild}|\text{yes}) = \frac{4}{9}$$

$$P(\text{humidity} = \text{normal}|\text{no}) = \frac{1}{5}$$

$$P(\text{humidity} = \text{normal}|\text{yes}) = \frac{6}{9}$$

$$P(\text{wind} = \text{normal}|\text{no}) = \frac{2}{5}$$

$$P(\text{wind} = \text{normal}|\text{yes}) = \frac{6}{9}$$

$$P(\text{no}|\text{sunny}, \text{mild}, \text{normal}, \text{normal}) = \frac{5}{14} \cdot \frac{3}{5} \cdot \frac{2}{5} \cdot \frac{1}{5} \cdot \frac{2}{5} / P(\text{sunny}, \text{mild}, \text{normal}, \text{normal})$$

$$P(\text{yes}|\text{sunny}, \text{mild}, \text{normal}, \text{normal}) = \frac{9}{14} \cdot \frac{2}{9} \cdot \frac{4}{9} \cdot \frac{6}{9} \cdot \frac{6}{9} / P(\text{sunny}, \text{mild}, \text{normal}, \text{normal})$$

$$P(\text{sunny}, \text{mild}, \text{normal}, \text{normal}) = \frac{5}{14} \cdot \frac{3}{5} \cdot \frac{2}{5} \cdot \frac{1}{5} \cdot \frac{2}{5} + \frac{9}{14} \cdot \frac{2}{9} \cdot \frac{4}{9} \cdot \frac{6}{9} \cdot \frac{6}{9}$$

## Atrybuty ciągłe

**Funkcja gęstości:** zastąpienie  $P(a_i = v_i | c = d)$  przez  $g_{d,i}(v_i)$ , gdzie  $g_{d,i}$  jest funkcją gęstości atrybutu  $a_i$  w klasie  $d$  – zakładany rozkład normalny z parametrami estymowanymi jako:

$$m_{T_{c=d}}(a_i) = \frac{1}{|T_{c=d}|} \sum_{x \in T_{c=d}} a_i(x)$$

$$s_{T_{c=d}}(a_i) = \sqrt{\frac{1}{|T_{c=d}| - 1} \sum_{x \in T_{c=d}} (a_i(x) - m_{T_{c=d}}(a_i))^2}$$

**Dyskretyzacja:** często lepsze, chociaż bardziej pracochłonne podejście.

## Obsługa brakujących wartości

**Tworzenie modelu:** pomijanie brakujących wartości przy estymacji prawdopodobieństw  $P(a_i = v_i | c = d)$ ,

**Predykcja:** pomijanie prawdopodobieństw  $P(a_i = v_i | c = d)$  jeśli wartość  $a_i$  nie jest znana dla klasyfikowanego przykładu.

# Właściwości naiwnego klasyfikatora bayesowskiego

- Prosty koncepcyjnie, implementacyjnie i obliczeniowo.
- Odporny na nadmierne dopasowanie.
- Niewymagający strojenia parametrów.
- Naruszenie założenia o niezależności nie wyklucza wartości predykcyjnej.
- Skuteczny jeśli:
  - trzeba uwzględnić nieznaczący wpływ znacznej liczby atrybutów (np. klasyfikacja tekstu),
  - liczba przykładów jest stosunkowo mała w porównaniu z liczbą atrybutów,
  - występują liczne brakujące wartości.