

Klasyfikacja w oparciu o przykłady



(ang. instance based learning)
Wykład 2, 14/10/2003

Plan wykładu

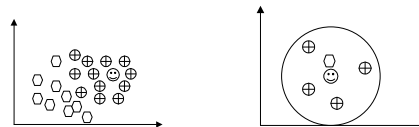
- Wprowadzenie
- Metoda k najbliższych sąsiadów
 - Miary podobieństwa
 - Praktyczne problemy
 - Redukcja zbędnych przykładów
 - Redukcja szumu w danych
 - Wyznaczanie wag atrybutów
 - Lokalna regresja
- Naiwna metoda wnioskowania Bayesowskiego
- Sieć Bayesowska

Lazy vs. eager learning

- **Eager learning model:**
 - Np. metody drzew decyzyjnych, reguł decyzyjnych, czy grupowania danych:
„Konstruuje się jasny opis funkcji docelowej na podstawie przykładów trenujących„
- **Lazy learning model:**
 - Np. klasyfikacja w oparciu o przykłady:
„Nie konstruuje się wcześniej opisu funkcji docelowej. Ta konstrukcja odbywa się w momencie klastrowania nowego obiektu”

Przykład

Dwuwymiarowy zbiór danych: każdy obiekt jest opisany dwoma atrybutami (x, y) .
Są dwie klasy \oplus lub \circ



Algorytm najbliższego sąsiada (algorytm 1-NN)

- **Parametr wejściowy:**
 - Zbiór obiektów $P = \{ \langle x_i, f(x_i) \rangle \}$, gdzie f – funkcja docelowa, np. opis klas decyzyjnych.
 - x_q - obiekt do klasyfikowania
- **Parametr wyjściowy:**
 - wartość $f(x_q)$ (np. klasa decyzyjna, do której należy x_q)

Algorytm najbliższego sąsiada (algorytm 1-NN)

- **Ogólny schemat:**
 - Krok 1:** Poszukaj obiektu x_n najbliższego x_q .
 - Krok 2:** Wyznacz wartość $f(x_q)$ na podstawie wartości $f(x_n)$
- **Zalety:** Prosty, szybki algorytm
- **Wady:** Nieodporny na szumy!!!

Algorytm k najbliższych sąsiadów (algorytm k -NN)

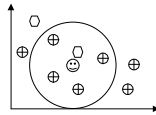
- **Ogólny schemat:**

Krok 1: Poszukaj k najbliższych obiektów (sąsiadów) dla x_q

Krok 2: Głosuj wśród k najbliższych sąsiadów w celu wyznaczenia klasy, do której należy x_q .

1-NN, decyzja jest \ominus

5-NN, decyzja jest \oplus



Zaleta: Bardziej odporny na szumy

Algorytm k -NN: Problemy

- Wyznaczanie miary podobieństwa (funkcji odległości) między obiektami.
- Głosowanie w celu wyznaczenia klasy, do której należy nowy obiekt.
- Wyznaczanie liczby k najbliższych obiektów potrzebnych dla klasyfikowania nowego obiektu.

Miary podobieństwa

- Niech każdy obiekt x będzie zdefiniowany wektorem wartości: $\langle a_1(x), \dots, a_n(x) \rangle$

- **Lokalna odległość :**

$$d(x, y) = |a_i(x) - a_i(y)|$$

- **Odległość między obiektami :**

$$\text{distance}(x, y) = F(d_1(x, y), \dots, d_n(x, y))$$

- **Odległość euklidesowa:**

$$\text{Eu-distance}(x, y) = [\sum_i (a_i(x) - a_i(y))^2]^{1/2}$$

- **Odległość miejska (Manhattan):**

$$\text{distance}(x, y) = \sum_i |a_i(x) - a_i(y)|$$

- **Miary podobieństwa:**

$$\text{sim}(x, y) = 1/(1 + \text{distance}(x, y))$$

Metody głosowania

- **Zasada większościowa:**

$$\hat{f}(x_q) \leftarrow \underset{v \in V}{\text{argmax}} \sum_{i=1}^k \delta(v, f(x_i))$$

$$\text{gdzie } \delta(x, y) = \begin{cases} 1 & \text{jesli } x = y \\ 0 & \text{wpp.} \end{cases}$$

- **Ważona odległość:**

$$\hat{f}(x) \leftarrow \underset{v \in V}{\text{argmax}} \sum_{i=1}^k w_i \delta(v, f(x_i)) \quad \text{gdzie } w_i = \frac{1}{1 + d(x_0, x_i)^2}$$

Wyznaczanie parametru k

- Jeśli k jest małe, algorytm nie jest odporny na szumy \rightarrow jakość klasyfikacji jest niska.
- Jeśli k jest duże, koszt obliczenia jest większy \rightarrow algorytm jest czasochłonny.
- Jak wybierać odpowiednią wartość k ?
- **Idea:**
 - Wykonuj test typu krosvalidacji dla kilku różnych wartości k .
 - Wybierz wartość k , która daje najwyższą jakość klasyfikacji.

Praktyczne problemy w algorytmie k -NN

1. Dane zawierają szumy:

Rozwiązanie: usuwanie szumów

2. Atrybuty w różnym stopniu są ważne

Rozwiązanie:
wyznaczanie wag dla atrybutów
(lub selekcja istotnych atrybutów)

3. Funkcja docelowa nie jest dyskretna (wartości są rzeczywiste)

Rozwiązanie:

- modyfikacja algorytmu k -NN
- lokalna regresja

Usuwanie szumu (noisy exemplars)



- **I strategia:**
 - wyznacz odpowiedni parametr k .
- **II strategia:**
 - Oceń jakość klasyfikacji każdego obiektu trenującego.
 - Usuń „słabe” obiekty.

Usuwanie szumu (c.d.)



- Dane: s_{min} , s_{max} – dolna i górna granica dokładności klasyfikacji.
- **Oceniać jakość klasyfikacji obiektu:**
 - Krok 1.** Wykonuj krosvalidację danych trenujących.
 - Krok 2.** Dla każdego obiektu x zanotuj procent obiektów dobrze klasyfikowanych s_x .
 - Krok 3.** Jeśli
 - $s_x < s_{min}$ to obiekt jest „słaby” (szum, trzeba usunąć)
 - $s_{min} < s_x < s_{max}$ obiekt jest „średnio dobry”

Wyznaczanie wag dla atrybutów

- **Motywacja:**
 - Niektóre atrybuty są ważniejsze niż inne
 - Niektóre atrybuty są ważne dla jednej klasy ale nie są ważne dla innej klasy.
- Jeśli w_1, w_2, \dots, w_n – wagi atrybutów to

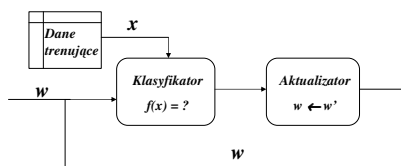
$$distance(x, y) = \sqrt{w_1^2 (a_1(x) - a_1(y))^2 + \dots + w_n^2 (a_n(x) - a_n(y))^2}$$

Wyznaczanie wag dla atrybutów (c.d.)

- **I model:** Każda klasa jest związana z jednym wektorem wag.
- **II model:** Jeden wspólny wektor wag dla wszystkich klas.
- **Idea:** Układ wag jest obliczony na podstawie wyniku klasyfikacji obiektów w zbiorze trenującym.
- **Algorytm sekwencyjnego poprawiania!**

Algorytm sekwencyjnego poprawiania

- **Ogólny schemat**
 - Krok 1.** zacznij od dowolnego wektora wag $w = [w_1, \dots, w_n]$.
 - Krok 2.** aktualizuj w , kiedy nowy obiekt trenujący x jest klasyfikowany.



Aktualizacja wektora wag

- **Wejście:** zbiór obiektów P
- **Wyjście:** wektor wag w
- Krok 1.** Dla $x \in P$ klasyfikuj x za pomocą pozostałych obiektów
- Krok 2.** Dla x , znajdź najbliższy obiekt y (wśród obiektów w zbiorze trenującym)
- Krok 3.** Dla każdego atrybutu a_i wyznacz $|a_i(x) - a_i(y)|$
- Krok 4.** Jeśli klasyfikacja jest prawidłowa to zwiększ wagę w_i , wpp. zmniejsz wagę w_i o Δw_i
(Δw_i jest odwrotnie proporcjonalne do $|a_i(x) - a_i(y)|$)
- Krok 5.** Jeśli wszystkie obiekty nie zostały testowane powróć do **Krok 1.**

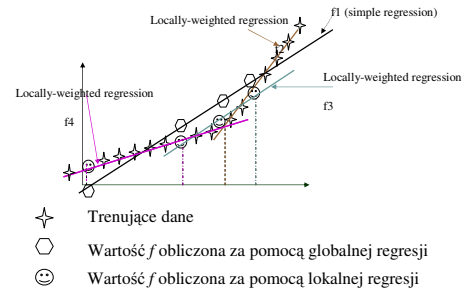
Wyznaczanie wartości rzeczywistej funkcji docelowej

- Modyfikacja algorytmu k -NN:

$$\hat{f}(x_q) \leftarrow \frac{\sum_{i=1}^k w_i f(x_i)}{\sum_{i=1}^k w_i} \quad \text{gdzie } w_i = \frac{1}{1 + d(x_q, x_i)^2}$$

- Lokalna regresja:
Znajdź lokalną aproksymację dla funkcji f , która najbardziej pasuje do obiektów w otoczeniu x_q .
- Liniowa regresja
- Kwadratowa regresja

Lokalna regresja



Metody Bayesowskie

- Naiwna metoda klasyfikacji Bayesowskiej
- Sieć Bayesowska
- Kombinacja z wiedzą dziedzinową

Podstawowa teoria

- Reguła Bayesowska: $P(h | D) = \frac{P(D | h)P(h)}{P(D)}$
gdzie
 - $P(h)$ = prawdopodobieństwo zajęcia hipotezy h
 - $P(D)$ = prawdopodobieństwo otrzymania zbioru treningowego D
 - $P(h|D)$ = prawdopodobieństwo h pod warunkiem, że D jest dany
 - $P(D|h)$ = prawdopodobieństwo D przy założeniu, że h zachodzi

Zasada Bayesowskiego uczenia się:

Szukanie najbardziej prawdopodobnej hipotezy mając zadany zbiór treningowy: (maksymalizacja hipotezy aposteriori h_{map})

$$\begin{aligned}
 h_{map} &= \max_{h \in H} P(h | D) \\
 &= \max_{h \in H} \frac{P(D | h)P(h)}{P(D)} \\
 &= \max_{h \in H} P(D | h)P(h)
 \end{aligned}$$

Podstawowe twierdzenia probabilistyczne

- Prawdopodobieństwo koniunkcji dwóch zdarzeń:

$$P(A, B) = P(A | B)P(B) = P(B | A)P(A)$$

- Prawdopodobieństwo sumy dwóch zdarzeń:

$$P(A + B) = P(A) + P(B) - P(AB)$$

- Wzór na prawdopodobieństwo całkowite:
jeśli zdarzenia A_1, \dots, A_n tworzą rozłączny podział przestrzeni probabilistycznej, to:

$$P(B) = \sum_{i=1}^n P(B | A_i)P(A_i)$$

Przykład

- Czy pacjent jest chory na raka?

Pacjent poddał testowi na obecność pewnego raka i dostał pozytywny wynik. Wynik testu jest

- prawidłowy (pozytywny) w 98% wśród chorujących na raka i
- prawidłowy (negatywny) w 97% wśród tych, którzy nie chorują na tego raka.
- Poza tym, 0.8% populacji choruje na badanego raka.

$$\begin{aligned}
 P(\text{cancer}) &= .008, P(\neg \text{cancer}) = .992 \\
 P(+ | \text{cancer}) &= .98, P(- | \text{cancer}) = .02 \\
 P(+ | \neg \text{cancer}) &= .03, P(- | \neg \text{cancer}) = .97 \\
 P(\text{cancer} | +) &= \frac{P(+ | \text{cancer})P(\text{cancer})}{P(+)} \\
 P(\neg \text{cancer} | +) &= \frac{P(+ | \neg \text{cancer})P(\neg \text{cancer})}{P(+)}
 \end{aligned}$$

Naiwna metoda Bayesa

- Załóżmy, że uczymy się funkcji celu $f: X \rightarrow V$, gdzie każdy obiekt jest opisany wektorem $\langle a_1, a_2, \dots, a_n \rangle$.
- Najbardziej prawdopodobna wartość $f(x)$ wynosi:

$$v = \max_{v_j \in V} P(v_j | a_1, a_2, \dots, a_n)$$

$$= \max_{v_j \in V} \frac{P(a_1, a_2, \dots, a_n | v_j) P(v_j)}{P(a_1, a_2, \dots, a_n)}$$

$$= \max_{v_j \in V} P(a_1, a_2, \dots, a_n | v_j) P(v_j)$$

- Naiwne założenie: „atrybuty są warunkowo niezależne”, tzn.

$$P(a_1, a_2, \dots, a_n | v_j) = \prod_i P(a_i | v_j)$$

Przykład: naiwna metoda

Zgadnij, czy odbędzie się gra w tenisa w dniu o warunkach pog.: $\langle \text{sunny, cool, high, strong} \rangle$ na podstawie nast. Zbioru danych:

Outlook	Temperature	Humidity	Windy	Class
sunny	hot	high	false	N
sunny	hot	high	true	N
overcast	hot	high	false	P
rain	mild	high	false	P
rain	cool	normal	false	P
rain	cool	normal	true	N
overcast	cool	normal	true	P
sunny	mild	high	false	N
sunny	cool	normal	false	P
rain	mild	normal	false	P
sunny	mild	normal	true	P
overcast	mild	high	true	P
overcast	hot	normal	false	P
rain	mild	high	true	N

$$p(y)p(\text{sunny}|y)p(\text{cool}|y)p(\text{high}|y)p(\text{strong}|y) = .005$$

$$p(n)p(\text{sunny}|n)p(\text{cool}|n)p(\text{high}|n)p(\text{strong}|n) = .021$$

przykład

Outlook	Temperature	Humidity	Windy	Class
sunny	hot	high	false	N
sunny	hot	high	true	N
overcast	hot	high	false	P
rain	mild	high	false	P
rain	cool	normal	false	P
rain	cool	normal	true	N
overcast	cool	normal	true	P
sunny	mild	high	false	N
sunny	cool	normal	false	P
rain	mild	normal	false	P
sunny	mild	normal	true	P
overcast	mild	high	true	P
overcast	hot	normal	false	P
rain	mild	high	true	N

$$P(p) = 9/14$$

$$P(n) = 5/14$$

outlook	
$P(\text{sunny} p) = 2/9$	$P(\text{sunny} n) = 3/5$
$P(\text{overcast} p) = 4/9$	$P(\text{overcast} n) = 0$
$P(\text{rain} p) = 3/9$	$P(\text{rain} n) = 2/5$
temperature	
$P(\text{hot} p) = 2/9$	$P(\text{hot} n) = 2/5$
$P(\text{mild} p) = 4/9$	$P(\text{mild} n) = 2/5$
$P(\text{cool} p) = 3/9$	$P(\text{cool} n) = 1/5$
humidity	
$P(\text{high} p) = 3/9$	$P(\text{high} n) = 4/5$
$P(\text{normal} p) = 6/9$	$P(\text{normal} n) = 2/5$
windy	
$P(\text{true} p) = 3/9$	$P(\text{true} n) = 3/5$
$P(\text{false} p) = 6/9$	$P(\text{false} n) = 2/5$

Algorytm „Naive Bayes”

1. Uczenie (zbiór przykładów)

- for (każda klasa decyzyjna v_j)
oszacuj $P(v_j)$
- for (każda wartość a_i na atrybucie a)
oszacuj $P(a_i | v_j)$

2. Klasyfikacja nowego obiektu(x)

$$v = \max_{v_j \in V} P(v_j) \prod_{a_i \in x} P(a_i | v_j)$$

typowe oszacowanie dla $P(a_i | v_j)$

$$P(a_i | v_j) \leftarrow \frac{n_c + mp}{n + m}$$

Gdzie:

- n : liczba przykładów z klasy v_j ; p : wstępne oszacowanie dla $P(a_i | v_j)$
- n_c : liczba przykładów z $a = a_i$; m : waga przekonań dla p

Sieć Bayesowska

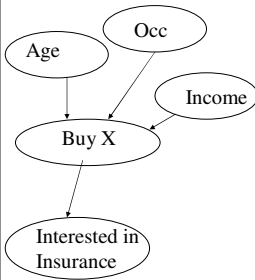
- Założenie „naive Bayes” jest zbyt ograniczone i prymitywne!
- Bez niego, obliczenia nie są wykonalne
- Sieć Bayesowska:
 - opisuje *warunkową niezależność między zbiorem atrybutów*, uwzględniając
 - *wiedzę ekspercką* o zależnościach między atrybutami i
 - *dane treningowe*.
 - DAG (direct acyclic graph)

Sieć Bayesowska

- Jest to acykliczny graf skierowany, gdzie
 - Wierzchołki: atrybuty
 - Krawędzie: zależność
 - Kierunki krawędzi: relacja przyczynowo-skutkowa
 - Do każdego atrybutu A, dołączona jest tablica prawdopodobieństw

$$P(A | B_1, \dots, B_n)$$
 gdzie B_1, \dots, B_n są bezpośrednimi poprzednikami atrybutu A w tym grafie

Przykład sieci



- Age, Occupation oraz Income decyduje, czy klient kupuje dany produkt.
- Jeśli klient kupuje produkt, to jego zainteresowanie ubezpieczeniem (interest in insurance) jest niezależne od Age, Occupation, Income.

$$P(\text{Age, Occ, Inc, Buy, Ins}) = P(\text{Age})P(\text{Occ})P(\text{Inc})P(\text{Buy}|\text{Age, Occ, Inc})P(\text{Int}|\text{Buy})$$

- **Stan wiedzy:** przy zadanej strukturze i warunkowych prawdopodobieństwach, istniejące algorytmy mogą wnioskować o atrybuty symboliczne i dyskretyzowane atrybuty ciągłe.

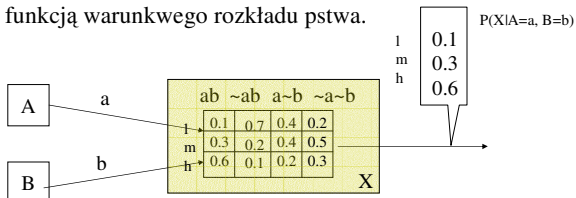
Podstawowy wzór:

$$P(x_1, \dots, x_n | M) = \prod_{i=1}^n P(x_i | Pa_i, M)$$

$$Pa_i = \text{parent}(x_i)$$

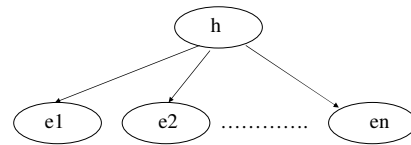
Wierzchołki jako funkcje

Każdy wierzchołek w sieci Bayesowskiej jest funkcją warunkowego rozkładu pstwa.



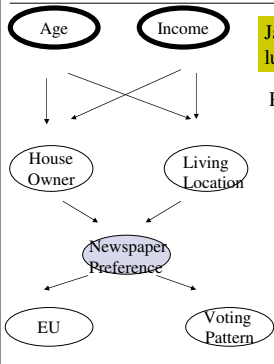
- wejście: wartość rodziców
- wyjście: rozkład pstwa własnych wartości

Przypadek szczególny: „naive Bayes”



$$P(e_1, e_2, \dots, e_n, h) = P(h) P(e_1 | h) \dots P(e_n | h)$$

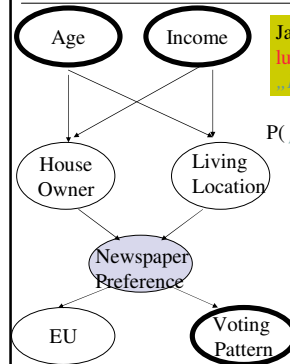
Wnioskowanie o sieci Bayesowskiej:



Jaka jest szansa, że **bogaci i starzy** ludzie kupują „Sun”?

$$P(\text{paper} = \text{Sun} | \text{Age} > 60, \text{Income} > 60k)$$

Wnioskowanie o sieci Bayesowskiej:

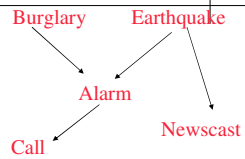


Jaka jest szansa, że **bogaci i starzy** ludzie głoszący na partię X kupują „Daily Mail”?

$$P(\text{paper} = \text{DM} | \text{Age} > 60, \text{Income} > 60k, \text{Vote} = X)$$

Uczenie Bayesowskie

B	E	A	C	N
~b	e	a	c	n
b	~e	~a	~c	n
.....				



Dane: pełna lub częściowa obserwacja przypadków

Szukane: parametry i struktura

Metody uczenia się:

EM (Expectation Maximisation)

-Uzupełnić brakujące dane za pomocą bieżącej aproksymacji parametrów;

-Aproksymować parametry za pomocą wypełnionych danych

Gradient Ascent Training

Gibbs Sampling (MCMC)

Bibliografia

- Lenz M., Bartsch-Sporl B., Burkhard H., Wess S. (1998). *Case-based reasoning technology. From foundations to applications*. Prinfen-Verlag Berlin Heidelberg. LNAI, Vol. 1400.
- Aha D. (1992). *Tolerating noisy, irrelevant, and novel attributes in instance-based learning algorithms*. International Journal of Man-Machine Studies 36(2): 267-287.