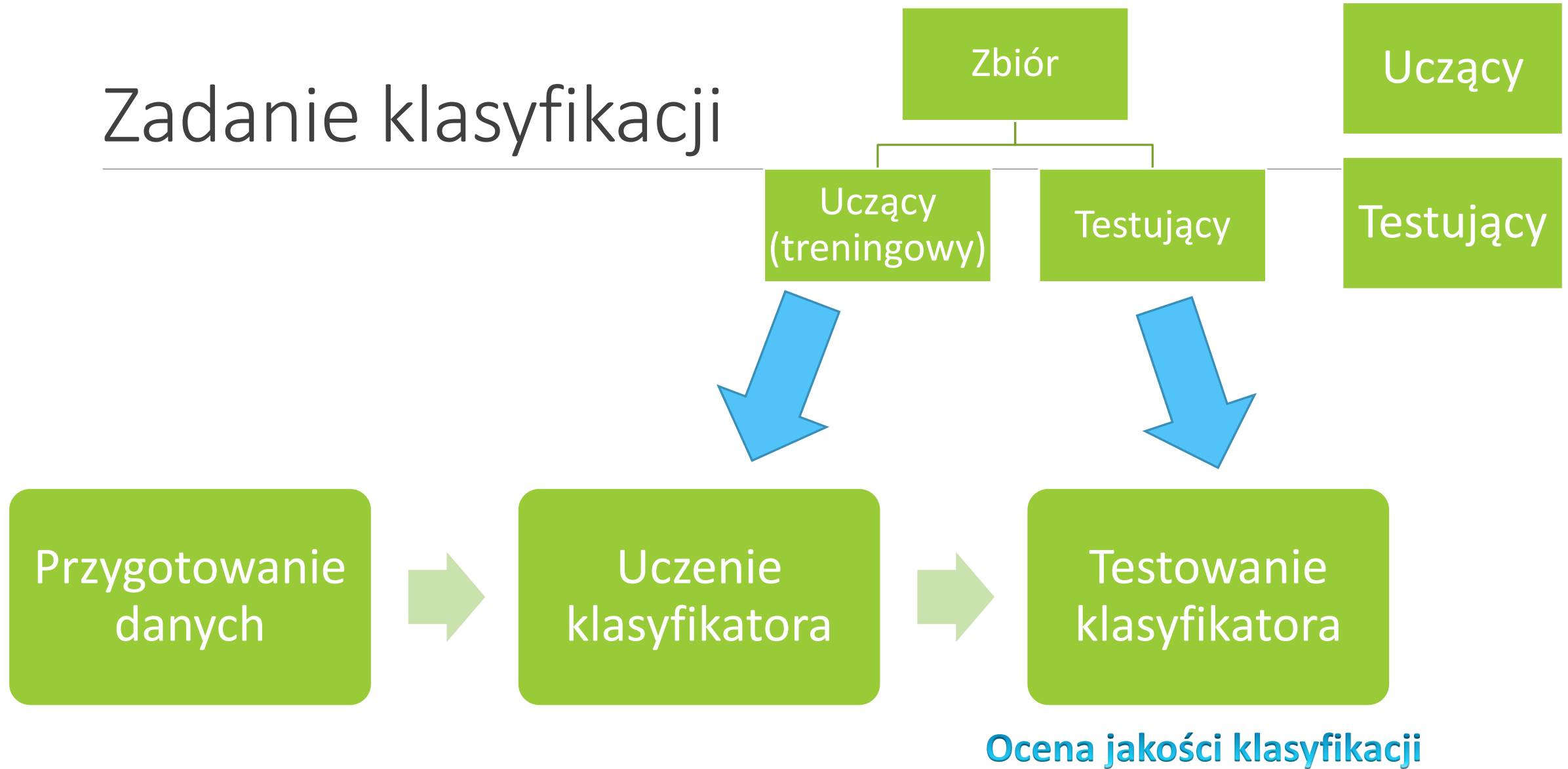


Uczenie maszynowe

ZADANIE KLASYFIKACJI CZ. 2

Zadanie klasyfikacji



k -krotna walidacja krzyżowa (k -fold *cross-validation*)

- Podział oryginalnego zbioru na k podzbiorów (losowy): Każdy z nich wykorzystywany jest jako zbiór testujący, a pozostałe $k-1$ podzbiorów jako zbiór uczący. Dla każdej z k iteracji wykonuje się analizę, a na końcu wyniki są uśredniane.
- Najczęściej k równe 5 lub 10



Naiwny klasyfikator Bayesa

- Klasyfikator statystyczny, oparty na twierdzeniu Bayesa:

Przykład X klasyfikujemy jako pochodzący z tej klasy C_i , dla której wartość $P(C_i|X)$, $i = 1, 2, \dots, m$, jest największa

gdzie X oznacza przykład z nieznaną klasą. Każdy przykład jest opisany jako wektor o n wymiarach $X = (x_1, \dots, x_n)$;

$P(C/X)$ prawdopodobieństwo a-posteriori, że przykład X należy do klasy C

- **Założenie: wzajemna niezależność atrybutów (naiwność ...)**

P-stwo warunkowe

- P-stwo warunkowe zajścia zdarzenia B pod warunkiem zajścia zdarzenia A:

$$P(\text{B} \mid \text{A}) = \frac{P(\text{A and B})}{P(\text{A})}$$

- Zadanie:

70% twoich koleżanek/kolegów lubi czekoladę, a 35% czekoladę i truskawki. Ile procent z tych, którzy lubią czekoladę, lubi również truskawki?

$$P(\text{truskawki} \mid \text{czekolada}) = \frac{P(\text{czekolada} \cap \text{truskawki})}{P(\text{czekolada})} = \frac{0,35}{0,7} = 0,5$$

Naiwny klasyfikator Bayesa

Wyznaczanie prawdopodobieństwa a-posteriori, że przykład X należy do klasy C na podstawie twierdzenia Bayesa:

$$P(C|X) = (P(X|C) * P(C))/P(X),$$

gdzie: **$P(C)$** oznacza prawdopodobieństwo a-priori wystąpienia klasy C (tzn. prawdopodobieństwo, że dowolny przykład należy do klasy C);

$P(X|C)$ oznacza prawdopodobieństwo a-posteriori, że X należy do klasy C ;

$P(X)$ oznacza prawdopodobieństwo a-priori wystąpienia przykładu X .

Wyprowadzenie:

$$P(C|X) = \frac{P(C \cap X)}{P(X)} = \frac{P(X|C) \times P(C)}{P(X)}, \text{ bo } P(X|C) = \frac{P(C \cap X)}{P(C)}$$

Naiwny klasyfikator Bayesa

- Dany jest zbiór treningowy D składający się z n przykładów
- Załóżmy, że atrybut decyzyjny przyjmuje m różnych wartości definiując m różnych klas C_i , $i = 1, \dots, m$
- Niech s_i oznacza liczbę przykładów z D należących do klasy C_i
- Klasyfikator Bayesa przypisuje nieznany przykład X do tej klasy C_i , dla której wartość $P(C_i|X)$ jest największa
- Prawdopodobieństwo $P(X)$ jest stałe dla wszystkich klas - klasa C_i , dla której wartość $P(C_i|X)$ jest największa, to klasa C_i , dla której wartość $P(X|C_i) * P(C_i)$ jest największa
- Wartości $P(C_i)$ zastępujemy estymatorami s_i/n (względna częstością klasy C_i), lub zakładamy, że wszystkie klasy mają to samo prawdopodobieństwo $P(C_1) = P(C_2) = \dots = P(C_m)$

$$P(C|X) = (P(X|C) * P(C))/P(X),$$

Naiwny klasyfikator Bayesa

- W jaki sposób obliczyć $P(X|C_i)$?
- Dla dużych zbiorów danych, o dużej liczbie deskryptorów, obliczenie $P(X|C_i)$ będzie bardzo kosztowne
- Przyjmujemy założenie o **niezależności atrybutów**
- Założenie o niezależności atrybutów prowadzi do następującej formuły:

$$P(X|C_i) = \prod_{j=1}^n P(x_j | C_i)$$

- Prawdopodobieństwa $P(x_1|C_i)$, $P(x_2|C_i)$, ..., $P(x_n|C_i)$ można estymować w oparciu o zbiór treningowy następująco:

jeżeli j-ty atrybut jest atrybutem kategoriowym, to $P(x_j|C_i)$ estymujemy względną częstością występowania przykładów z klasy C_i posiadających wartość x_j dla j-tego atrybutu, (s_{ij}/s_i)

jeżeli j-ty atrybut jest atrybutem ciągłym, to $P(x_j|C_i)$ estymujemy **funkcją gęstości Gaussa**

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

(zakładając rozkład normalny wartości atrybutów)

$$P(C|X) = (P(X|C) * P(C))/P(X),$$

Zbiór czyKupiKomputer

■ Atrybuty:

- wiek { ≤ 30 , 31...40, > 40 }
- dochód {duży, średni, mały}
- czy student? {tak, nie}
- zdolność kredytowa {umiark, doskonała}
- czy kupi komputer? {tak, nie} → KLASA

wiek	dochod	stud	zdolKred	czyKupi
≤ 30	duży	nie	umiark	nie
≤ 30	duży	nie	doskonala	nie
31 ... 40	duży	nie	umiark	tak
> 40	sredni	nie	umiark	tak
> 40	mały	tak	umiark	tak
> 40	mały	tak	doskonala	nie
31 ... 40	mały	tak	doskonala	tak
≤ 30	sredni	nie	umiark	nie
≤ 30	mały	tak	umiark	tak
> 40	sredni	tak	umiark	tak
≤ 30	sredni	tak	doskonala	tak
31 ... 40	sredni	nie	doskonala	tak
31 ... 40	duży	tak	umiark	tak
> 40	sredni	nie	doskonala	nie

Zadanie czyKupiKomputer

Nowy przypadek: **X = (wiek = ≤30, dochód = duży, student = tak, zdolność kredytowa = doskonała)**

$$P(\text{czyKupi} = \text{tak}) = P(C1) = 9/14 \approx 0,643$$

$$P(\text{czyKupi} = \text{nie}) = P(C2) = 5/14 \approx 0,357$$

$$P(\text{wiek} = \leq 30 | \text{czyKupi} = \text{tak}) = 2/9 \approx 0,222$$

$$P(\text{wiek} = \leq 30 | \text{czyKupi} = \text{nie}) = 3/5 \approx 0,6$$

$$P(\text{dochód} = \text{duży} | \text{czyKupi} = \text{tak}) = 2/9 \approx 0,222$$

$$P(\text{dochód} = \text{duży} | \text{czyKupi} = \text{nie}) = 2/5 \approx 0,4$$

$$P(\text{student} = \text{tak} | \text{czyKupi} = \text{tak}) = 6/9 \approx 0,667$$

$$P(\text{student} = \text{tak} | \text{czyKupi} = \text{nie}) = 1/5 \approx 0,2$$

$$P(\text{zdolKred} = \text{doskonała} | \text{czyKupi} = \text{tak}) = 3/9 \approx 0,333$$

$$P(\text{zdolKred} = \text{doskonała} | \text{czyKupi} = \text{nie}) = 3/5 \approx 0,6$$

wiek	dochod	stud	zdolKred	czyKupi
≤ 30	duży	nie	umiark	nie
≤ 30	duży	nie	doskonała	nie
31 ... 40	duży	nie	umiark	tak
> 40	sredni	nie	umiark	tak
> 40	mały	tak	umiark	tak
> 40	mały	tak	doskonała	nie
31 ... 40	mały	tak	doskonała	tak
≤ 30	sredni	nie	umiark	nie
≤ 30	mały	tak	umiark	tak
> 40	sredni	tak	umiark	tak
≤ 30	sredni	tak	doskonała	tak
31 ... 40	sredni	nie	doskonała	tak
31 ... 40	duży	tak	umiark	tak
> 40	sredni	nie	doskonała	nie

Zadanie czyKupiKomputer

$$\begin{aligned} P(X|\text{czyKupi} = \text{tak}) &= P(\text{wiek} = \leq 30|\text{czyKupi} = \text{tak}) \times P(\text{dochód} = \text{duży}|\text{czyKupi} = \text{tak}) \times \\ &\times P(\text{student} = \text{tak}|\text{czyKupi} = \text{tak}) \times P(\text{zdolKred} = \text{doskonała}|\text{czyKupi} = \text{tak}) = \\ &= 0,222 \times 0,222 \times 0,667 \times 0,333 \approx 0,011 \end{aligned}$$

$$P(X|\text{czyKupi} = \text{nie}) = 0,6 \times 0,4 \times 0,2 \times 0,6 \approx 0,029$$

$$P(X|\text{czyKupi} = \text{tak}) \times P(\text{czyKupi} = \text{tak}) = 0,011 \times 0,643 \approx 0,007$$

$$P(X|\text{czyKupi} = \text{nie}) \times P(\text{czyKupi} = \text{nie}) = 0,029 \times 0,357 \approx 0,010$$

Przypadkowi X zostanie przypisana etykieta klasy:

$$\text{czyKupi} = \text{nie}$$

Problem zmiennych numerycznych

- Gra w golfa w zależności od pogody
- Atrybuty:
 - Temperatura
 - Wilgotność
 - czyGramy? Play {yes, no}

**Jak policzyć p-stwo wystąpienia
poszczególnych wartości
np. temperatury 69?**

Rozwiązanie – dyskretyzacja!

Temperature	Humidity	Play
85	85	no
80	90	no
65	70	no
72	95	no
71	80	no
83	78	yes
70	96	yes
68	80	yes
64	65	yes
69	70	yes
75	80	yes
75	70	yes
72	90	yes
81	75	yes

Zadanie 1

Na podstawie zbioru czyKupiKomputer zaklasyfikuj dwa następujące przypadki:

- 1) **X = (wiek = 31 ... 40, dochód = sredni, student = nie, zdolność kredytowa = umiark)**
- 2) **swój własny przypadek**

wiek	dochod	stud	zdolKred	czyKupi
≤ 30	duży	nie	umiark	nie
≤ 30	duży	nie	doskonala	nie
31 ... 40	duży	nie	umiark	tak
> 40	sredni	nie	umiark	tak
> 40	mały	tak	umiark	tak
> 40	mały	tak	doskonala	nie
31 ... 40	mały	tak	doskonala	tak
≤ 30	sredni	nie	umiark	nie
≤ 30	mały	tak	umiark	tak
> 40	sredni	tak	umiark	tak
≤ 30	sredni	tak	doskonala	tak
31 ... 40	sredni	nie	doskonala	tak
31 ... 40	duży	tak	umiark	tak
> 40	sredni	nie	doskonala	nie