

# Zadanie klasyfikacji - drzewa decyzyjne, wykł. 5

Joanna Jędrzejowicz

Instytut Informatyki

# Klasyfikacja danych - przykłady

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

- Bankowość- czy dane podanie o udzielenie kredytu jest obarczone dużym czy małym ryzykiem, czy dana transakcja kartą jest oszustwem,
- Edukacja - umieszczenie studenta w odpowiedniej grupie z uwzględnieniem jego potrzeb,
- Medycyna - diagnozowanie, czy występuje dana choroba,
- Telekomunikacja - czy dany klient zrezygnuje z usług firmy telekomunikacyjnej (zbiór churn)
- Informatyka - czy wiadomosc pocztowa jest spamem

W zadaniu **klasyfikacji** jest ustalona zmienna celu (inaczej: klasa) np. w powyższych przykładach: ryzyko duże - tak lub nie, grupa studencka o numerach 1, ..., 5, diagnoza - jest choroba- tak lub nie itd.

# Na czym polega zadanie klasyfikacji danych

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Niech  $C$  oznacza zbiór klas. Zakładamy, że algorytm uczący się klasyfikacji otrzymuje, jako wejście, zbiór danych treningowych (uczących)

$$LD = \{ \langle d, c \rangle \mid d \in D, c \in C \} \subset D \times C,$$

gdzie  $D$  jest zbiorem rekordów (wierszy atrybutów)  
 $d = (w_1^d, \dots, w_n^d)$  gdzie  $w_i^d$  dla  $i = 1, \dots, n$  jest wartością  $i$ -tego atrybutu,  $n$  - jest liczbą atrybutów.

# Na czym polega zadanie klasyfikacji danych, cd

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Algorytm uczenia klasyfikatora służy do znalezienia najlepszej możliwej aproksymacji  $\bar{f}$  nieznanej funkcji  $f$  takiej, że  $f(d) = c$ .

Następnie  $\bar{f}$  może być użyta do znalezienia klasy  $\bar{c} = \bar{f}(\bar{d})$  dla dowolnego  $\bar{d}$  takiego, że  $(\bar{d}, \bar{c}) \notin LD$ , zatem algorytm będzie mógł być użyty do znajdowania klasy dla przykładów, których nie użyto w uczeniu.

Zbiór danych  $LD$  składa się z dwóch podzbiorów  $LD = TD \cup TS$ , gdzie  $TD$  - jest zbiorem **treningowym**,  $TS$  - zbiorem **testowym**.

# Dokładność klasyfikacji

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Jak liczyć dokładność klasyfikatora:

- budujemy klasyfikator w oparciu o zbiór treningowy - dla każdego rekordu określona jest jego klasa,
- używając zbudowanego klasyfikatora, dla każdego wiersza ze zbioru testowego wyznaczamy klasę oraz porównujemy z rzeczywistą klasą,
- dokładność określamy jako stosunek liczby poprawnie sklasyfikowanych wierszy ze zbioru testowego do mocy tego zbioru.

# Ewaluacja- cross-validation - sprawdzian krzyżowy

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Ewaluacja metodą '10-fold cross-validation' polega na dokonaniu podziału zbioru danych na 10 części  $X_1, \dots, X_{10}$  i wykonaniu 10 doświadczeń.

W  $i$ -tym doświadczeniu  $X_i$  jest zbiorem testowym, a pozostałe  $X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_{10}$  podzbiorów stanowi zbiór treningowy; w każdym doświadczeniu oblicza się dokładność klasyfikatora, a następnie liczy się średnią dokładność.

# Na czym polega zadanie klasyfikacji danych, cd

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Klasyfikacja jest nazywana **uczeniem nadzorowanym**, gdyż algorytm uczący ma w zbiorze treningowym pełne dane - informacje, do jakiej klasy należy każdy wektor danych.

**Przewidywanie** (predykcja) jest podobnym zadaniem do klasyfikacji- w przewidywaniu wynik dotyczy przyszłości, np. przewidywanie cen akcji po upływie pół roku, przewidywanie zwycięzcy rozgrywek sportowych na podstawie wyników z przeszłości itd. Metody i techniki wykorzystywane do klasyfikacji mogą być również użyte, przy poczynieniu dodatkowych założeń, do przewidywania.

# Przykładowy zbiór danych 'czy kupi komputer?'

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

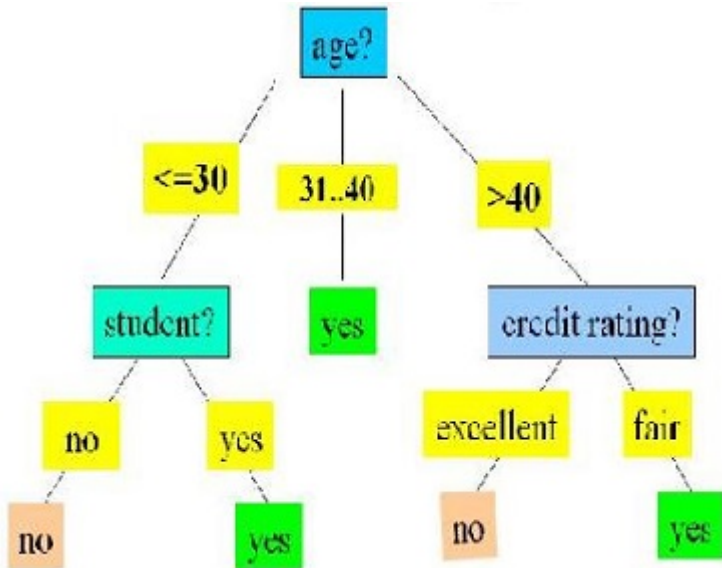
wiek	dochod	stud	zdolKred	czyKupi
$\leq 30$	duży	nie	umiark	nie
$\leq 30$	duży	nie	doskonała	nie
31...40	duży	nie	umiark	tak
$> 40$	sredni	nie	umiark	tak
$> 40$	mały	tak	umiark	tak
$> 40$	mały	tak	doskonała	nie
31...40	mały	tak	doskonała	tak
$\leq 30$	sredni	nie	umiark	nie
$\leq 30$	mały	tak	umiark	tak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



# Drzewo decyzyjne

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz



# Jak tworzy się drzewo decyzyjne?

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

- Z korzeniem drzewa związany jest pełny zbiór danych
- Każdy wierzchołek różny od liścia jest węzłem decyzyjnym
- Z każdym węzłem decyzyjnym związany jest atrybut
- Zbiór danych jest dzielony zgodnie z wartościami tego atrybutu i gałęzie prowadzą do następnych węzłów odpowiadających możliwym wartościom atrybutu
- Jeżeli z węzłem jest związany zbiór danych z jednej klasy, to węzeł jest liściem etykietowanym tą klasą
- Jeżeli w węźle występują dane z takimi samymi wartościami wszystkich atrybutów (nie można dalej dzielić!), to przyjmuje się głosowanie większościowe

# Kiedy można zastosować drzewo decyzyjne

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

- Zbiór treningowy powinien być bogaty i różnorodny, zapewniający reprezentatywną grupę typów rekordów, których klasyfikacja może być potrzebna w przyszłości,
- klasy zmiennej celu muszą być dyskretne, to znaczy nie można zastosować analizy drzew decyzyjnych do ciągłej zmiennej celu,
- drzewa decyzyjne starają się stworzyć zbiór liści, które są 'najczystsze', to znaczy gdy każdy z rekordów w danym liściu należy do tej samej klasy.

Uwaga: jeżeli wartości atrybutu są rzeczywiste, to trzeba podzielić zbiór wartości na pewną liczbę przedziałów

# Jakie mogą się pojawić problemy

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Co zrobić, jeżeli dany węzeł zawiera różne rekordy, z niejednakową wartością zmiennej celu i nie można dokonać dalszego podziału?

id	wiek	dochod	czyKupi
20	>40	duzy	tak
21	>40	duzy	tak
22	>40	duzy	nie
23	>40	duzy	tak

# Jak dokonywać wyboru atrybutów w poszczególnych węzłach- algorytm C4.5

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Dla każdego atrybutu liczy się **zysk informacji** związany z tym atrybutem i wybiera się atrybut, dla którego wartość jest największa.

Założmy, że (w ustalonym węźle) zbiór danych  $S$  składa się z  $s$  wierszy. Założmy, że zbiór klas (wartości zmiennej celu) składa się z  $m$  różnych elementów  $C_1, \dots, C_m$ . Niech  $s_i$  oznacza liczbę wierszy z klasy  $C_i$ .

Ilość informacji niezbędna do zaklasyfikowania rekordu z danego zbioru:

$$I(s_1, \dots, s_m) = - \sum_{i=1}^m p_i \times \log p_i$$

gdzie

$$p_i = \frac{s_i}{s}$$

# Jak dokonywać wyboru atrybutów w poszczególnych węzłach, cd

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Jeżeli  $A$  jest atrybutem, który przyjmuje  $v$  możliwych wartości  $a_1, \dots, a_v$ , to atrybut  $A$  dzieli zbiór  $S$  na  $v$  podzbiorów  $S_1, \dots, S_v$ . Niech  $s_{ij}$  oznacza liczbę wierszy z klasy  $C_i$  należących do podzbioru  $S_j$  (czyli wartość atrybutu  $A$  jest w tym podzbiorze równa  $a_j$ ).

Entropia

$$E(A) = \sum_{j=1}^v \frac{s_{1j} + \dots + s_{mj}}{s} \times I(s_{1j}, \dots, s_{mj})$$

$$I(s_{1j}, \dots, s_{mj}) = - \sum_{i=1}^m p_{ij} \times \log p_{ij}$$

gdzie

$$p_{ij} = \frac{s_{ij}}{|S_j|}$$

# Jak dokonywać wyboru atrybutów w poszczególnych węzłach, cd

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Zysk informacji dla atrybutu  $A$ :

$$Gain(A) = I(s_1, \dots, s_m) - E(A)$$

# Obliczenia dla przykładowych danych

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Mamy 4 atrybuty:

- wiek (z trzema możliwymi wartościami: ' $\leq 30$ ', ' $31..40$ ', ' $> 40$ '),
- dochód (z trzema możliwymi wartościami)
- student (tak lub nie)
- zdolKred (dwie wartości)

Przyjmijmy klasa  $C_1 = \text{'tak'}$ , klasa  $C_2 = \text{'nie'}$



# Wyliczenie zysku informacji dla atrybutu 'wiek'

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

$$I(s_1, s_2) = -\frac{9}{14} \log \frac{9}{14} - \frac{5}{14} \log \frac{5}{14} = 0,940$$

# Zysk informacji dla atrybutu 'wiek'

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Dla wartości ' $\leq 30$ ':

$$s_{11} = 2, s_{21} = 3, I(s_{11}, s_{21}) = -\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5} = 0,971$$

Dla wartości ' $31..40$ ':

$$s_{12} = 4, s_{22} = 0, I(s_{12}, s_{22}) = 0$$

Dla wartości ' $\geq 40$ ':

$$s_{13} = 3, s_{23} = 2, I(s_{13}, s_{23}) = -\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} = 0,971$$

## Zysk informacji dla atrybutu 'wiek', cd

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

$$E(wiek) = \frac{5}{14} \times I(s_{11}, s_{21}) + \frac{4}{14} \times I(s_{12}, s_{22}) + \frac{5}{14} \times I(s_{13}, s_{23}) = 0,69$$

$$Gain(wiek) = I(s_1, s_2) - E(wiek) = 0,940 - 0,694 = 0,246$$

**Zadanie** Policzyc zysk informacji dla pozostałych atrybutów.  
Sprawdzić, że największa wartość jest osiągnięta dla atrybutu  
'wiek'. Narysować drzewo decyzyjne ([wyklad5.xls](#))

# Błędy popełniane przez modele klasyfikacyjne

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

- błąd treningowy klasyfikatora - stosunek niepoprawnie zaklasyfikowanych rekordów zbioru **treningowego** do łącznej liczby rekordów w zbiorze treningowym,
- błąd testowy klasyfikatora - jak wyżej dla zbior **testowego**,
- błąd generalizacji - oczekiwany błąd na zbiorze nowych rekordów, zakłada się że minimalizacja błędu testowego prowadzi do minimalizacji błędu generalizacji

# Przykładowy zbiór testowy

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

wiek	dochod	stud	zdolKred	czyKupi
$\leq 30$	mały	tak	umiark	tak
$> 40$	duży	nie	doskonała	tak
$31 \dots 40$	duży	tak	umiark	tak

Jaka jest błąd testowy, jaki treningowy?

# Przeuczenie klasyfikatora

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Problem silnego dopasowania drzewa decyzyjnego (por. przykład):

- błąd treningowy mały,
- dane mogą być mało reprezentatywne,
- może powodować duży błąd generalizacji

# Przycinanie drzewa decyzyjnego (ang. pruning)

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Usunięcie poddrzewa i zastąpienie lisciem, któremu przypisuje się etykietą dominującą w zbiorze rekordów związanych z wierzchołkiem.

Przykład. Odcinamy wierzchołek **zdolKred** i zastępujemy lisciem 'tak'. Jak zmieni się błąd treningowy i testowy?

# Przycinanie drzewa decyzyjnego

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Metody wykorzystujące strategię **przycinania wstępnego**: na etapie konstrukcji drzewa zatrzymuje się tworzenie nowych węzłów korzystając z nast. kryteriów:

- miara jakości podziału zbioru treningowego (np. zysk informacji) jest poniżej zadanego progu,
- zbiór treningowy związany z danym wierzchołkiem nie jest dostatecznie liczny,
- rozkład klas (np. większość rekordów należy do jednej klasy)



# Przycinanie drzewa decyzyjnego, cd

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Metody wykorzystujące strategię **przycinania końcowego**:  
konstruuje się pełne drzewo, a następnie metodą wstępującą,  
poczynając od lisci i przesuwając się w kierunku korzenia  
dokonuje się przycięcia. Możliwe kryteria:

- wylicza się błąd klasyfikacji dla drzewa zredukowanego i porównuje z miarą dla drzewa niezredukowanego, jeżeli przycięcie pogarsza to przywraca się zredukowane poddrzewo
- ewentualnie kryteria jak dla przycinania wstępnego

# Reguły decyzyjne

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Korzystając z drzewa decyzyjnego możemy wygenerować **reguły decyzyjne** - wystarczy zapisać informację towarzyszącą przejściu od korzenia do liścia, np.

**if** wiek  $< 40 \wedge$  wiek  $\geq 31$  **then** czyKupi=**tak**

Dla reguł określa się dwie miary: **wsparcie** (ang. support) oraz **ufność** (ang. confidence).

Dla reguły **if A then B**

- **wsparcie**: stosunek liczby rekordów zawierających  $A$  i  $B$  do liczby wszystkich rekordów,
- **ufność**: stosunek liczby rekordów zawierających  $A$  i  $B$  do liczby rekordów zawierających  $A$ .
- Jakiek jest wsparcie i ufność dla powyższej reguły?

wsparcie  $\frac{4}{14}$ , ufność  $\frac{4}{4}$

# Metoda klasyfikacji Random Forest (Breiman 2001)

Zadanie  
klasyfikacji -  
drzewa  
decyzyjne,  
wykł. 5

Joanna  
Jędrzejowicz

Klasyfikator **zespołowy** wykorzystujący drzewa decyzyjne oraz następujące dwie zasady:

- 'bagging' - do budowy kolejnych drzew losuje się dane ze zbioru treningowego,
- przy generowaniu konkretnego drzewa, dla każdego wierzchołka losuje się  $M$  atrybutów i spośród nich wybiera jeden z największym zyskiem informacji ( $M$  jest parametrem).
- po skonstruowaniu drzew klasyfikuje się nowe przykłady stosując głosowanie większościowe.