

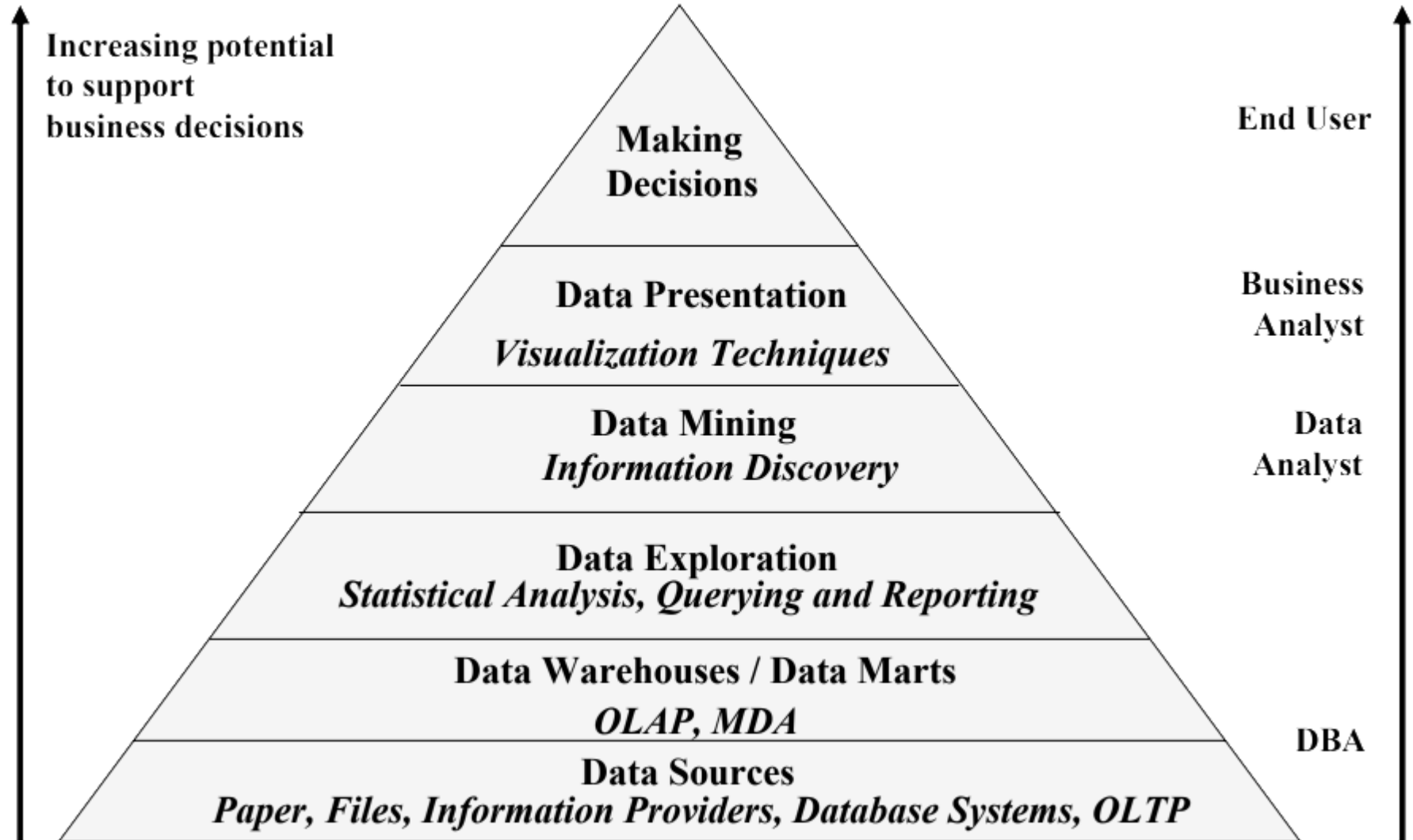
Indukcja reguł

w wykładzie wykorzystano:

1. materiały dydaktyczne przygotowane w ramach projektu *Opracowanie programów nauczania na odległość na kierunku studiów wyższych – Informatyka*
<http://wazniak.mimuw.edu.pl>
2. *Internetowy Podręcznik Statystyki*
<http://www.statsoft.pl/textbook/stathome.html>
3. J. Stefanowski – wykłady
4. Berthold, Borgelt, Höppner, Klawonn, *Guide to Intelligent Data Analysis*, Springer-Verlag London Limited 2010
5. Witten, Frank, *Data Mining Practical Machine Learning Tools and Techniques – WEKA*, Elsevier, San Francisco, 2005
6. *Data Science Central* - online resource for data practitioners

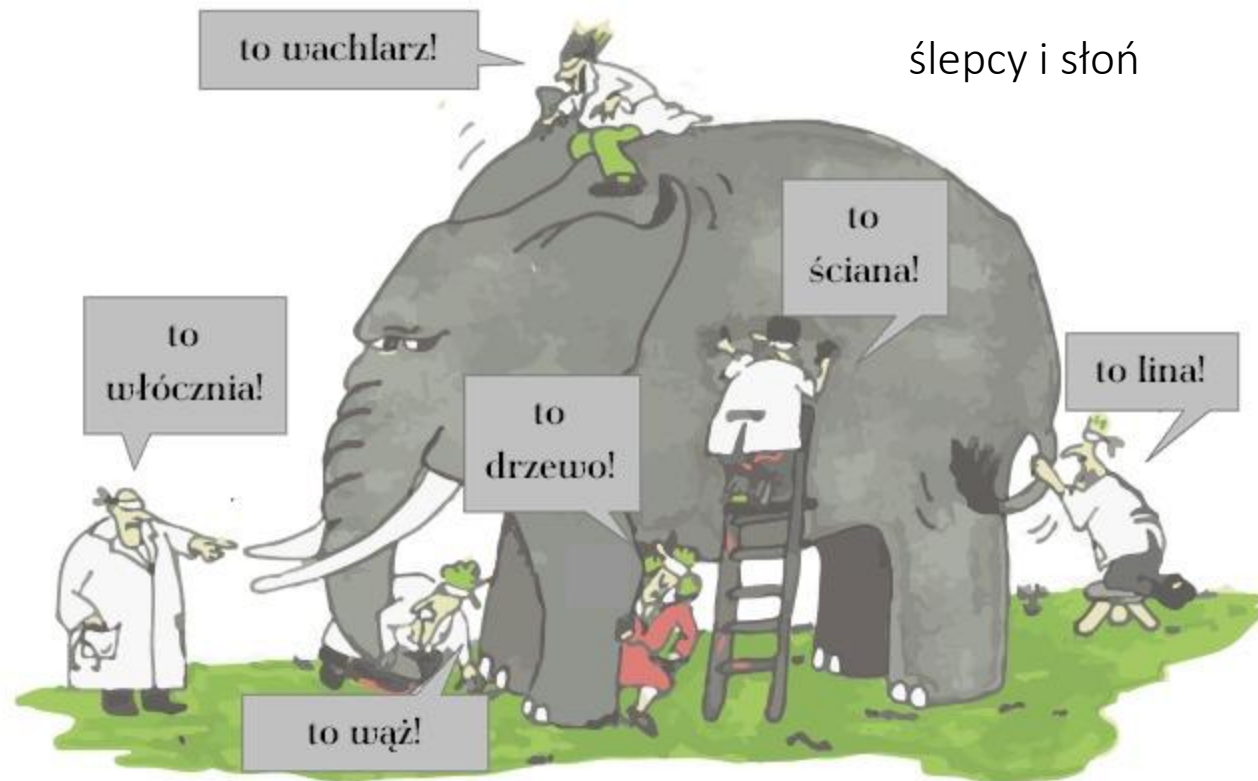
Krzysztof Regulski, WIMiIP, KISiM,
regulski@metal.agh.edu.pl

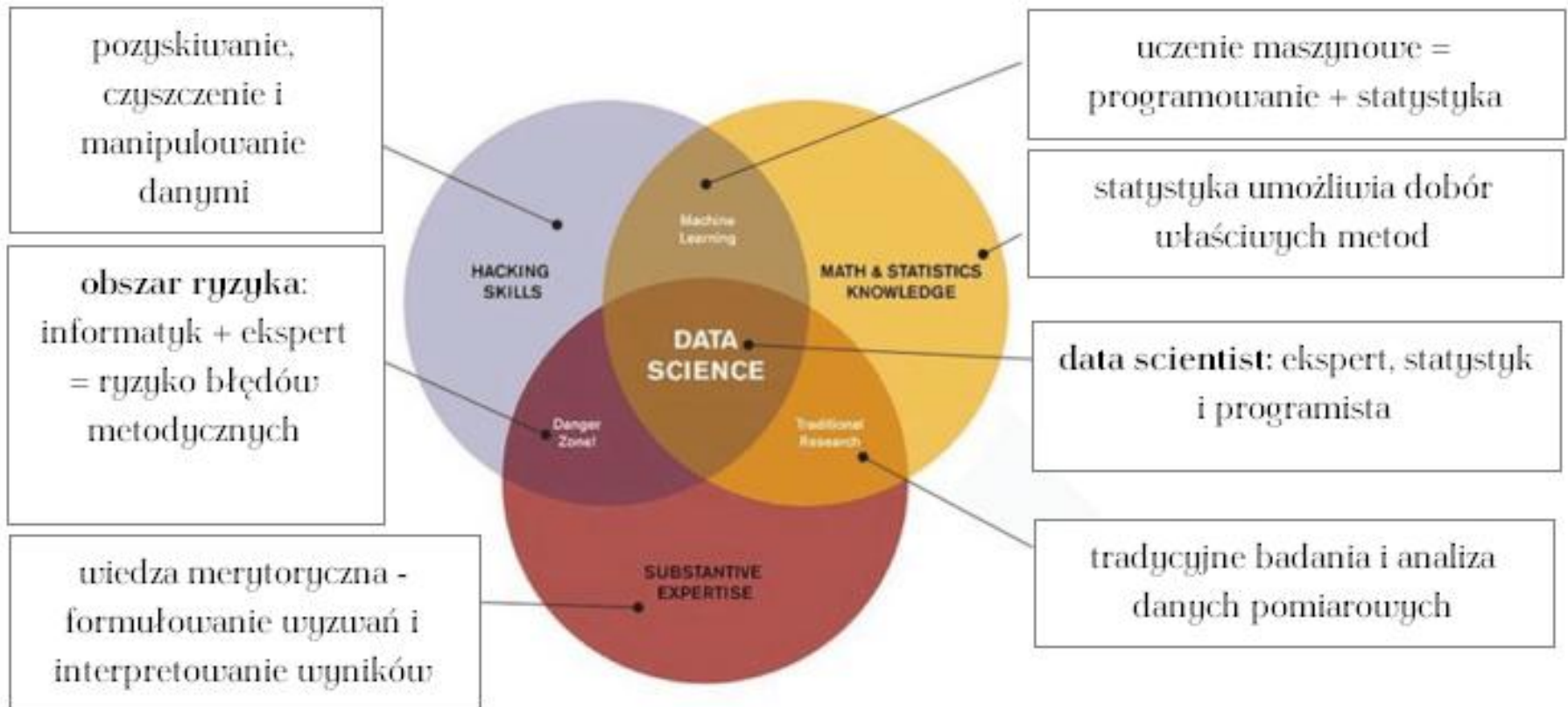
Data Mining and Business Intelligence



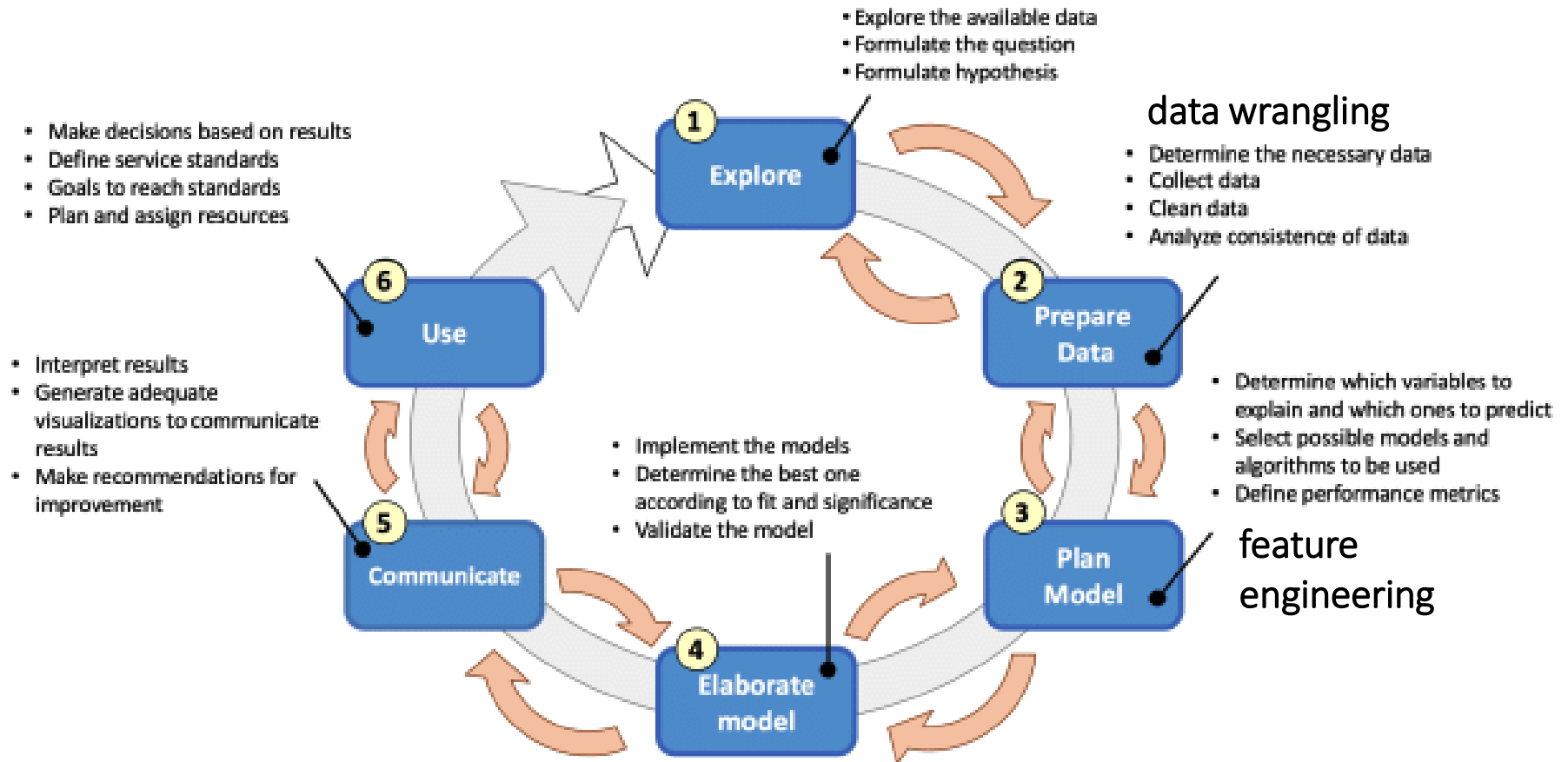
data science: *“a set of fundamental principles that support and guide the principled extraction of information and knowledge from data.”*

„data rich but information poor”





Knowledge Discovery from Data(KDD)

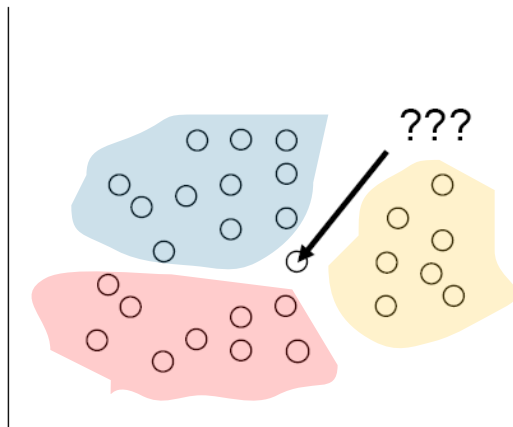


- **klasyfikacja:** przewidywanie wartości klasy na podstawie opisu
- **predykcja:** przewidywanie wartości ciągłej
- **analiza skupień:** poszukiwanie podobnych grup obiektów (skupień)
- **analiza asocjacji:** częste współwystępowanie sygnałów
- **analiza sekwencji:** analiza sekwencji sygnałów (obiektów) – analiza szeregów czasowych zmiennej jakościowej
- **charakterystyka:** tworzenie opisów grup
- **wizualizacja:** metody graficzne prezentacji wzorców
- **odkrywanie anomalii:** wykrywanie istotnych zmian (np. oszustw)

przede wszystkim prostota...

- proste algorytmy często wystarczają
- jest wiele przypadków prostych struktur danych:
 - jeden atrybut załatwia sprawę
 - wszystkie atrybuty mają podobny wpływ i są niezależne
 - logiczna struktura i mała liczba atrybutów odpowiednich dla drzew
 - zbiór kilku reguł logicznych
 - zależności pomiędzy grupami atrybutów
 - liniowa kombinacja atrybutów
 - wyraźne sąsiedztwo obiektów mierzone na podstawie odległości
 - wyraźne skupienia obiektów dla danych nieskategoryzowanych
 - zbiór obiektów dających się agregować
- skuteczność metod zależy od dziedziny/problemu badawczego

- Klasyfikacja jest metodą analizy danych, której celem jest **predykcja wartości** określonego **atrybutu** w oparciu o pewien zbiór **danych treningowych**.
- Obejmuje metody odkrywania **modeli** (tak zwanych **klasyfikatorów**) lub **funkcji** opisujących zależności pomiędzy **charakterystyką** obiektów a ich **zadaną klasyfikacją**.
- Odkryte modele klasyfikacji są wykorzystywane do klasyfikacji nowych obiektów o **nieznanej klasyfikacji**.



Wiele technik:

- statystyka,
- drzewa decyzyjne,
- sieci neuronowe, etc.

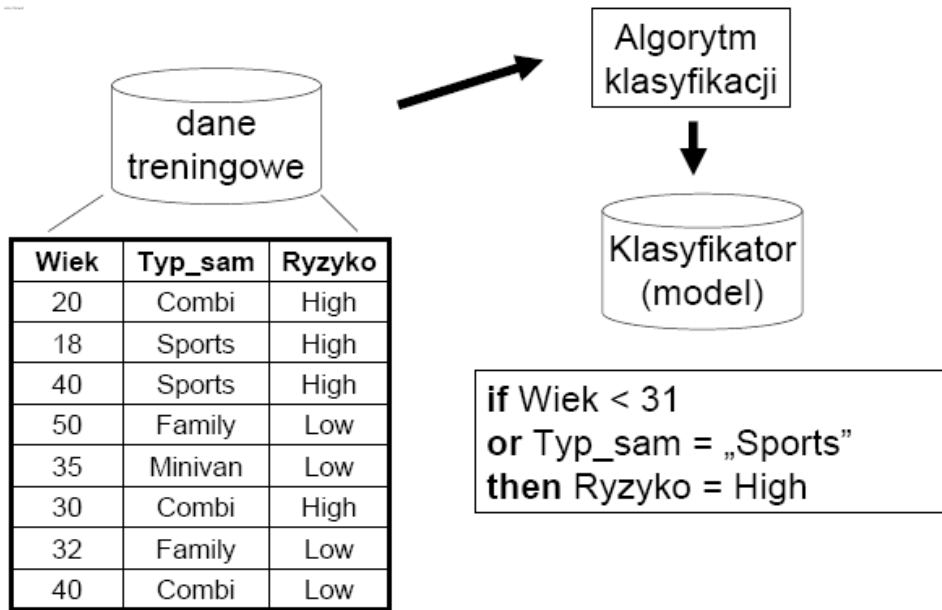
– Dane wejściowe

treningowy zbiór obserwacji będących listą wartości atrybutów opisowych (tzw. deskryptorów) i wybranego atrybutu decyzyjnego (*ang. class label attribute*)

– Dane wyjściowe

model (klasyfikator), **przydziela każdej krotce wartość atrybutu decyzyjnego** na podstawie wartości pozostałych atrybutów (deskryptorów, predyktorów)

Klasyfikacja – algorytm



Atrybut *Ryzyko* związany z informacją, że dany kierowca spowodował wcześniej wypadki czy nie spowodował wcześniej wypadku.

Jeżeli jest sprawcą kilku wypadków wartość atrybutu *Ryzyko* przyjmuje wartość High, w przypadku gdy nie spowodował żadnego wypadku atrybut *Ryzyko* przyjmuje wartość Low.

Atrybut *Ryzyko* jest **atrybutem decyzyjnym**.

W naszym przykładzie wynikiem działania algorytmu klasyfikacji jest klasyfikator w postaci pojedynczej reguły decyzyjnej: „Jeżeli wiek kierowcy jest mniejszy niż 31 lub typ samochodu sportowy to *Ryzyko* jest wysokie”.

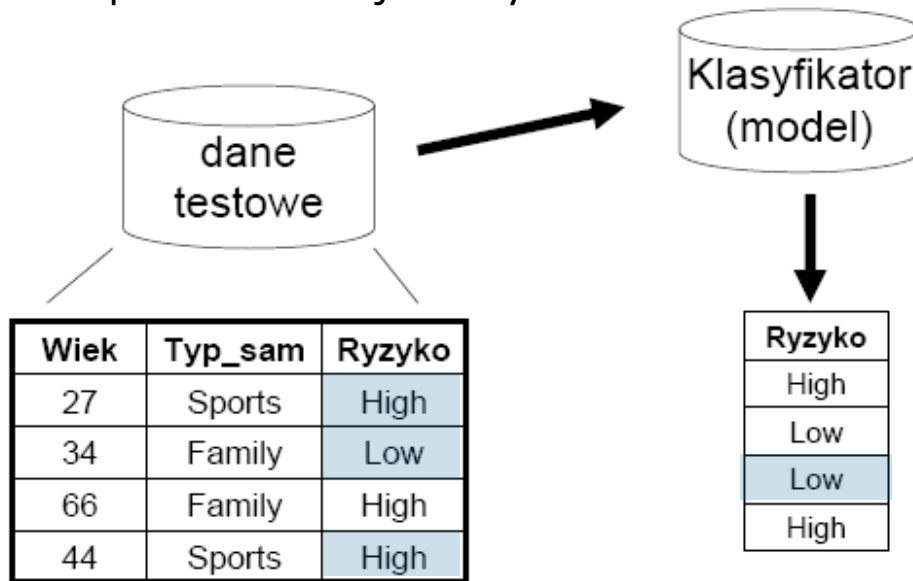
- Wynik klasyfikacji:
 - » Reguły klasyfikacyjne postaci *IF - THEN*
 - » Drzewa decyzyjne
- Istotną sprawą z punktu widzenia poprawności i efektywności modelu jest tzw. **dokładność modelu**.
- dla przykładów **testowych**, dla których **znane** są wartości atrybutu decyzyjnego, wartości te są porównywane z wartościami atrybutu decyzyjnego **generowanymi** dla tych przykładów przez klasyfikator.
- Miarą, która weryfikuje poprawność modelu jest **współczynnik dokładności**.

Współczynnik dokładności (*ang. accuracy rate*) = %
procent przykładów testowych
poprawnie zaklasyfikowanych przez model

Klasyfikacja – wynik

Jeżeli dokładność klasyfikatora jest **akceptowalna**, wówczas możemy wykorzystać klasyfikator do klasyfikacji **nowych danych**.

Celem klasyfikacji, jak pamiętamy jest przyporządkowanie nowych danych dla których wartość atrybutu decyzyjnego nie jest znana do odpowiedniej klasy.



Dokładność = $3/4 = 75\%$

Duży zbiór danych

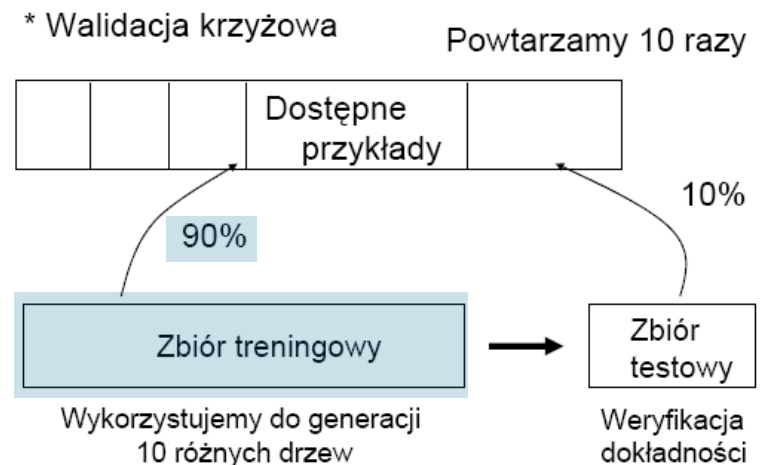
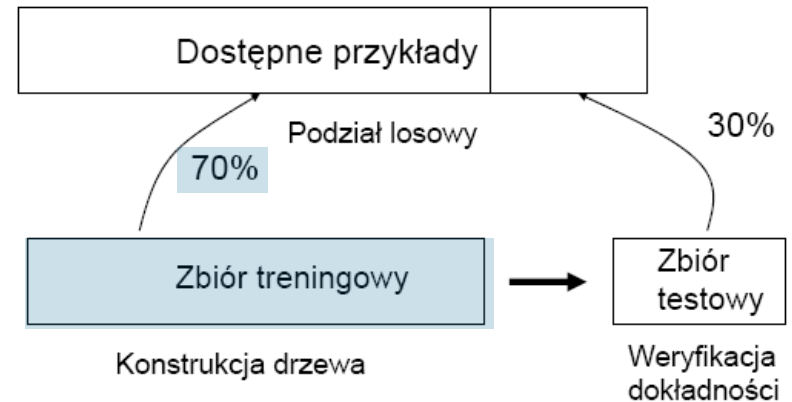
Mały zbiór danych

W przypadku zbioru przykładów o małej liczności stosujemy najczęściej metodę k-krotnej walidacji krzyżowej (tzw. krosvalidacji).

Początkowy zbiór przykładów jest losowo dzielony na k możliwie równych, wzajemnie niezależnych części S_1, S_2, \dots, S_k .

Zbiór treningowy stanowi $k-1$ części, k -ta część stanowi zbiór testowy. Sam klasyfikator konstruujemy k -krotnie. W ten sposób otrzymujemy k -klasyfikatorów

Po wybraniu klasyfikatora, klasyfikator konstruuje się raz jeszcze w oparciu o cały dostępny zbiór przykładów



Sprawdzian krzyżowy

cross-validation

1. Dzieli się dane na v rozłącznych części (wybranych losowo, losowanie bez zwracania).
2. Dla ustalonego wstępnie K wykonuje się analizę, by znaleźć **predykcję dla v -tej grupy** danych (używając pozostałych $v-1$ części danych jako przypadków "przykładowych").
3. Liczymy błąd predykcji. W przypadku regresji obliczamy sumę kwadratów reszt, przy klasyfikacji obliczamy dokładność, czyli procent przypadków zaklasyfikowanych poprawnie.
4. Na końcu v cykli **uśredniamy błędy**, otrzymując miarę jakości modelu.

Powyższe powtarzamy dla **różnych K** , wybierając jako najlepsze to K , dla którego otrzymujemy **najlepszą jakość** modelu.

wielokrotne repróbkiwanie

bootstrap

- metoda szacowania dokładności klasyfikatora dla **mało licznego zbioru** przykładów
- wykorzystuje losowanie ze zwracaniem
- tworzymy nowy zbiór losując n razy (z n -elementowego zbioru)
- niektóre przykłady będą się powtarzać w zbiorze treningowym, a inne przykłady w tym zbiorze nie wystąpią (dokładnie **0,368%** przykładów nie zostanie wylosowanych)
- nowy zbiór obejmie 63,2% przypadków (*0.632 bootstrap*)
- niewylosowane przykłady utworzą zbiór testowy, który wykorzystujemy do oceny dokładności otrzymanego klasyfikatora.

$$\left(1 - \frac{1}{n}\right)^n \approx e^{-1} = 0.368$$

- **klasyfikacja:** przewidywanie wartości klasy na podstawie opisu (wartości innych zmiennych)
- **predykcja:** przewidywanie wartości ciągłej, modelowanie funkcji ciągłych

Jeśli atrybut decyzyjny jest **ciągły** (numeryczny), mówimy o problemie **predykcji / regresji**.

Predykcja jest bardzo podobna do klasyfikacji. Jednakże celem predykcji jest zamodelowanie **funkcji ciągłej**, która by odwzorowywała **wartości** atrybutu decyzyjnego (**regresja**)

Rodzaje modeli klasyfikacyjnych:

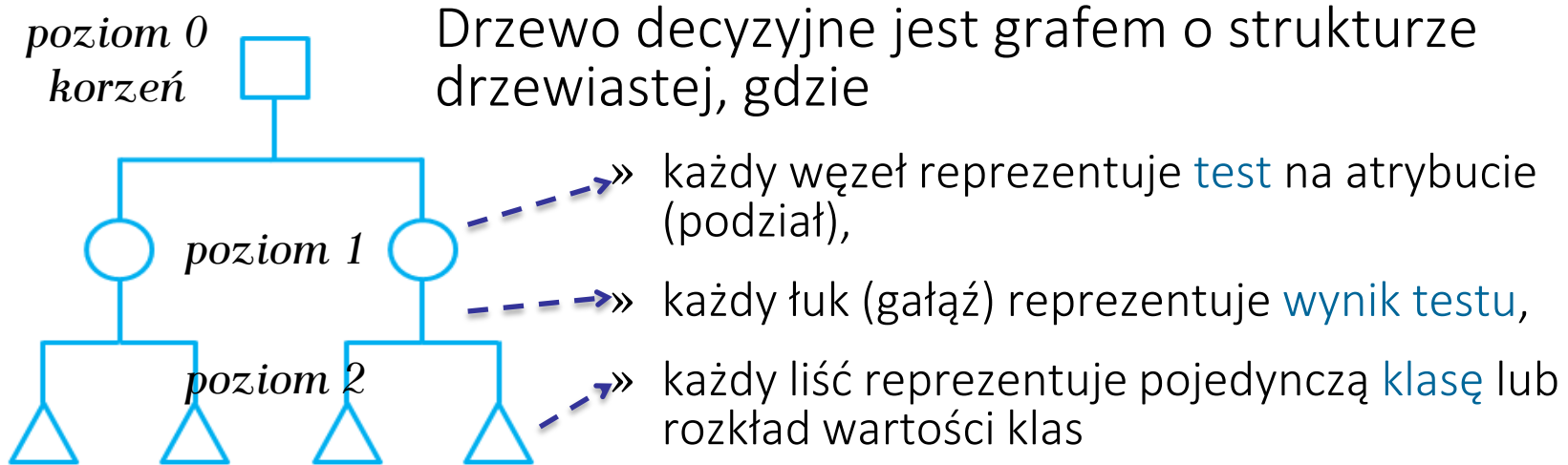
- » k-NN - k-najbliższe sąsiedztwo
(wartościowanie leniwe, lazy evaluation)
- » Klasyfikatory Bayes'owskie
(wartościowanie zachłanne, eager evaluation) - generative model approaches
- » Klasyfikacja poprzez indukcję drzew decyzyjnych
- » Klasyfikatory liniowe (LDA, FLA, etc.)
- » Discriminative Modelling Approaches
(Linear Classifier, Logistics Regression, SVM)
- » SVM – (Support Vector Machine) - Metoda wektorów nośnych
- » indukcja reguł, zbiory przybliżone, reguły asocjacyjne
- » Sieci Neuronowe, neuro-fuzzy
- » Analiza statystyczna, Metaheurystyki
(np. algorytmy genetyczne)
- » i inne...

Indukcja drzew klasyfikacyjnych

zmienna zależna: jakościowa

- rodzina algorytmów uczenia maszynowego
- uczenie nadzorowane (*supervised*), z nauczycielem
- dane są opisane kategoriami (etykietowane, *labelled*)
- uczenie wykorzystywane m.in.
 - w autonomicznych pojazdach do rozpoznawania obiektów,
 - w medycynie,
 - w ocenie wartości klientów,
 - chatterbotach,
 - itd.

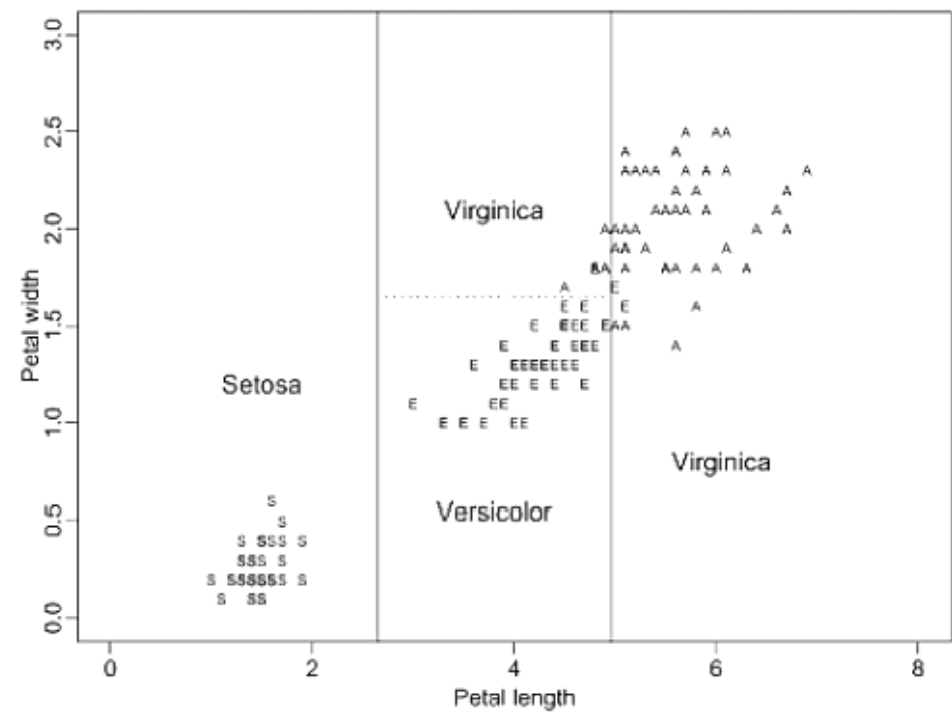
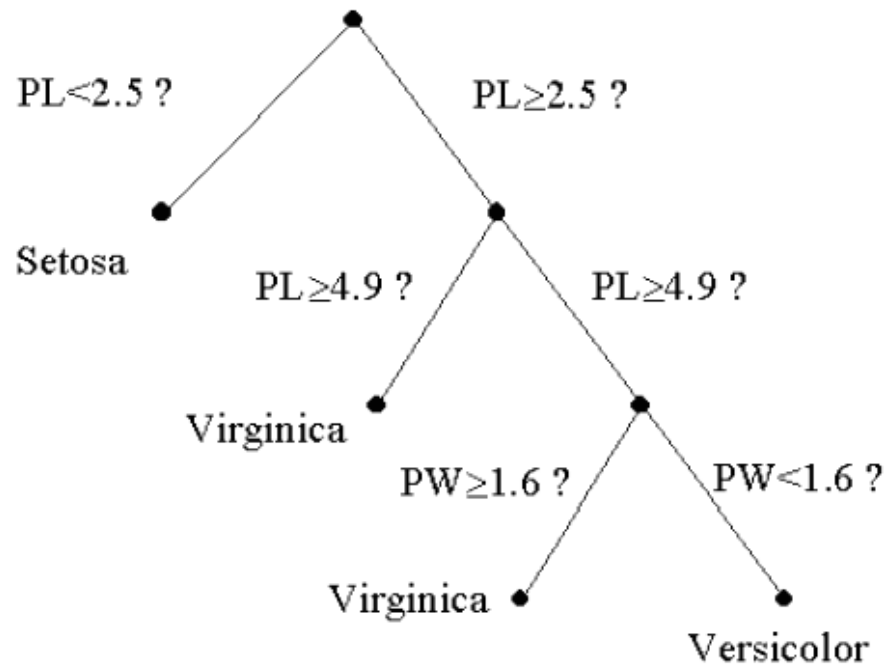
Klasyfikacja poprzez indukcję drzew decyzyjnych

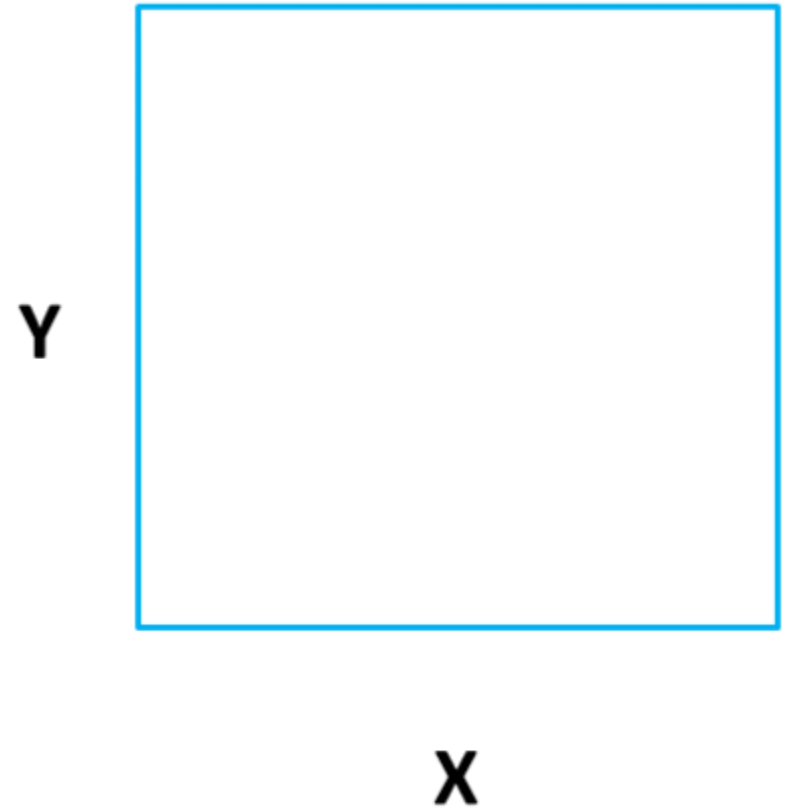


Drzewo decyzyjne dzieli zbiór treningowy na **partycje** do momentu, w którym każda partycja zawiera dane należące do **jednej klasy**, lub, gdy w ramach partycji dominują dane należące do jednej klasy

Każdy wierzchołek wewnętrzny (węzeł) drzewa zawiera tzw. **punkt podziału** (*ang. split point*), którym jest test na atrybucie (atrybutach), który dzieli zbiór danych na partycje

- *divide-and-conquer*
 - 1: wybierz atrybut podziału korzenia (pełnego zbioru)
utwórz gałąź dla każdej wartości atrybutu
 - 2: podziel zbiór na partycje
dla każdej gałęzi
 - 3: powtarzaj rekursywnie dla każdej gałęzi
używaj tylko instancji należących do partycji
- przerwij, jeśli wszystkie instancje należą do jednej klasy

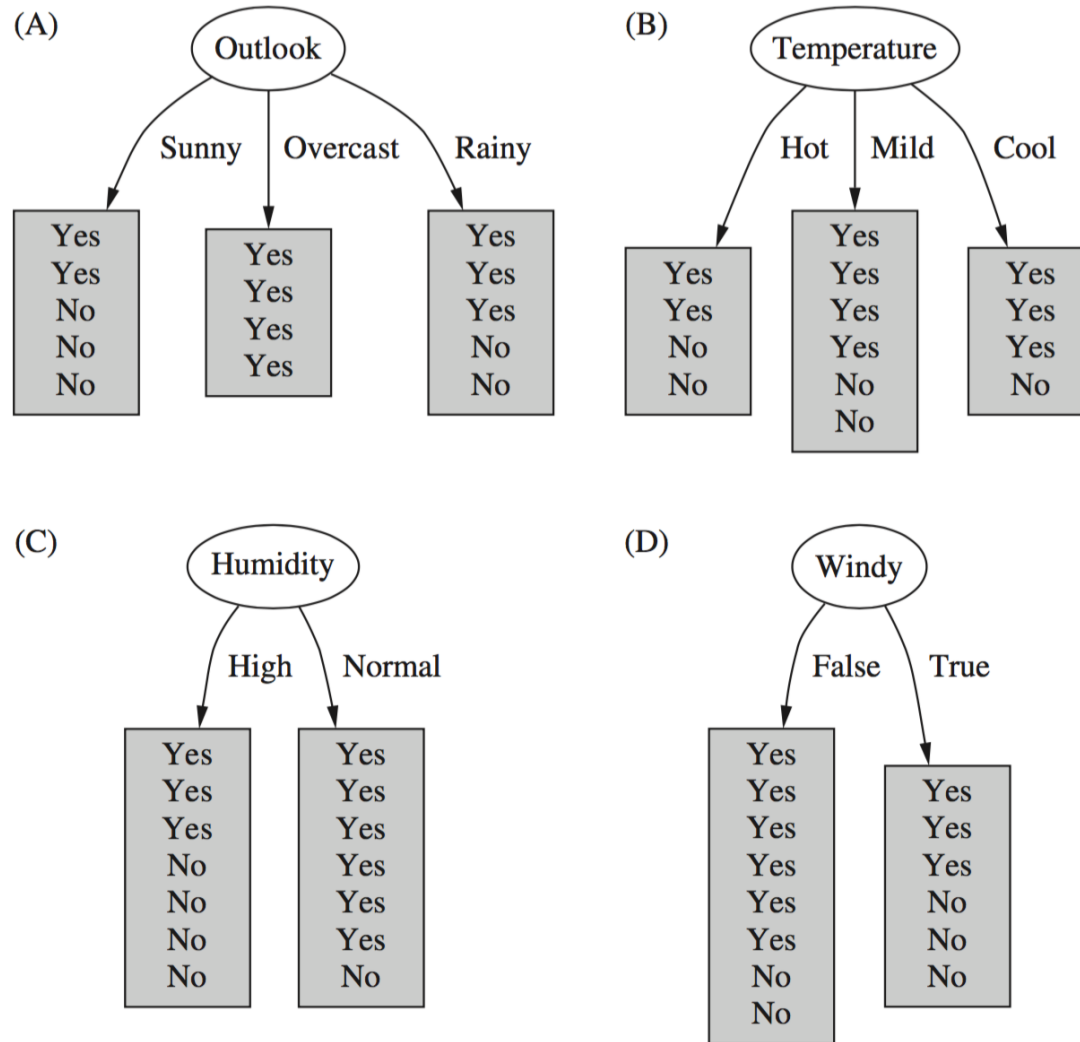




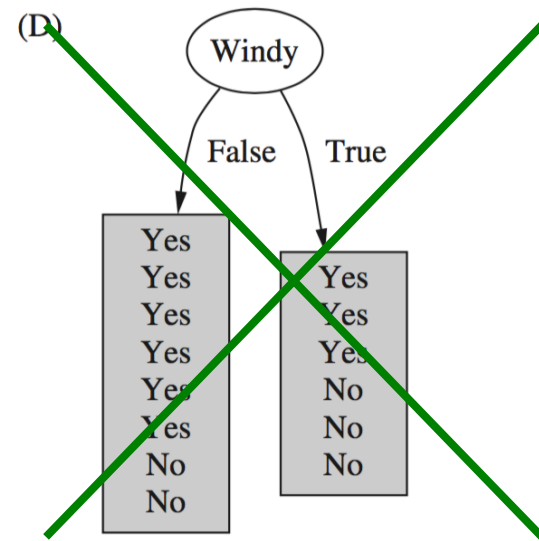
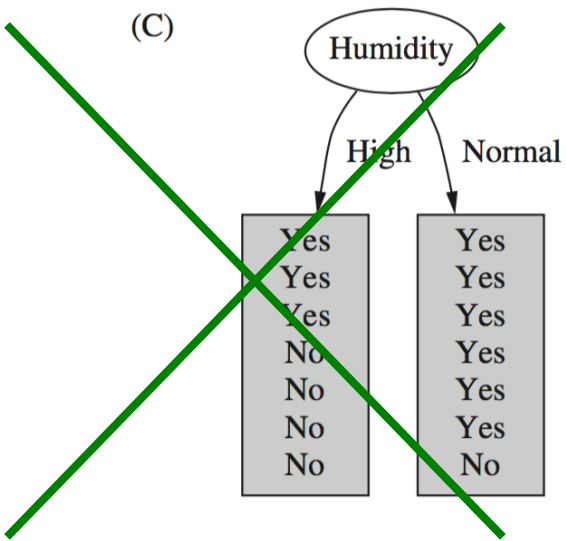
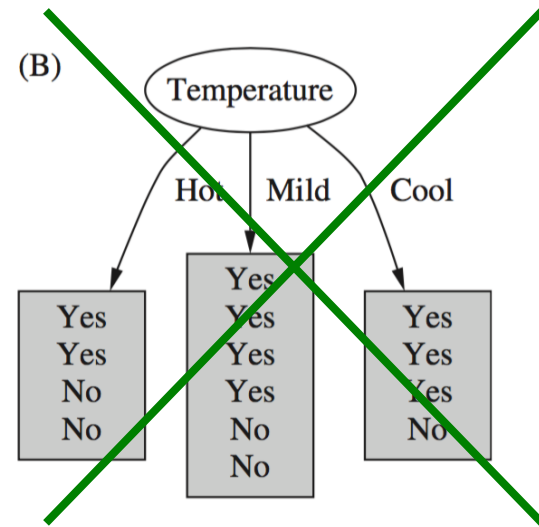
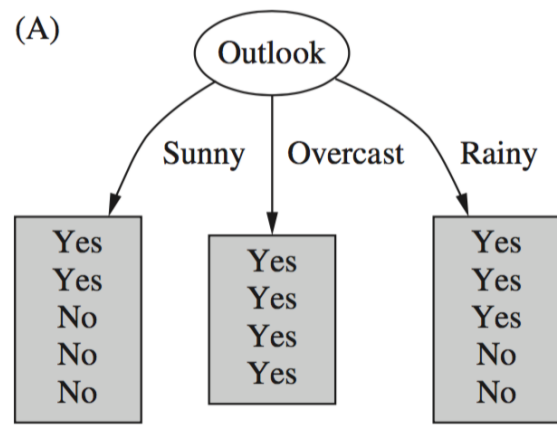
Czy dziś grać w golfa?

zachmurzenie	temperatura	wilgotność	wiatr	decyzja
słońce	gorąco	wysoka	słaby	nie
słońce	gorąco	wysoka	silny	nie
pochmurno	gorąco	wysoka	słaby	tak
deszcz	średnio	wysoka	słaby	tak
deszcz	chłodno	normalna	słaby	tak
deszcz	chłodno	normalna	silny	nie
pochmurno	chłodno	normalna	silny	tak
słońce	średnio	wysoka	słaby	nie
słońce	chłodno	normalna	słaby	tak
deszcz	średnio	normalna	słaby	tak
słońce	średnio	normalna	silny	tak
pochmurno	średnio	wysoka	silny	tak
pochmurno	gorąco	normalna	słaby	tak
deszcz	średnio	wysoka	silny	nie

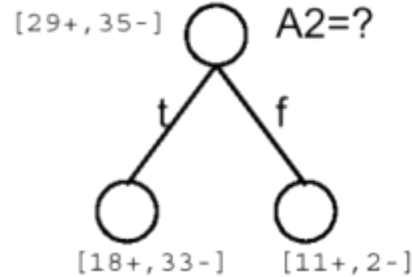
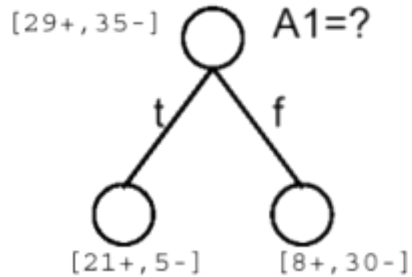
wybór predyktora (1)



wybór predyktora (2)



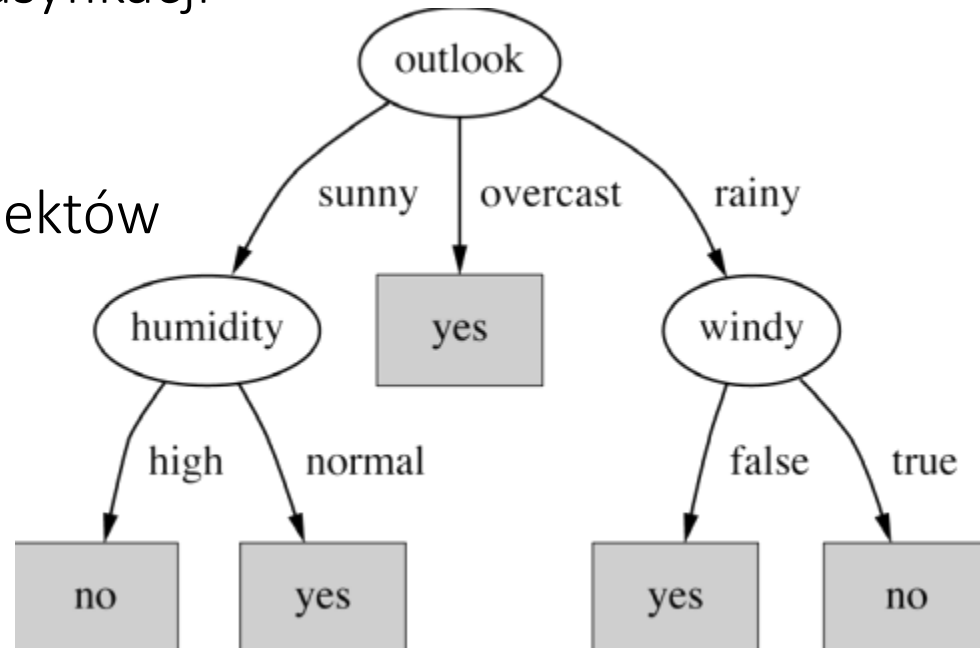
wybór predyktora (3)



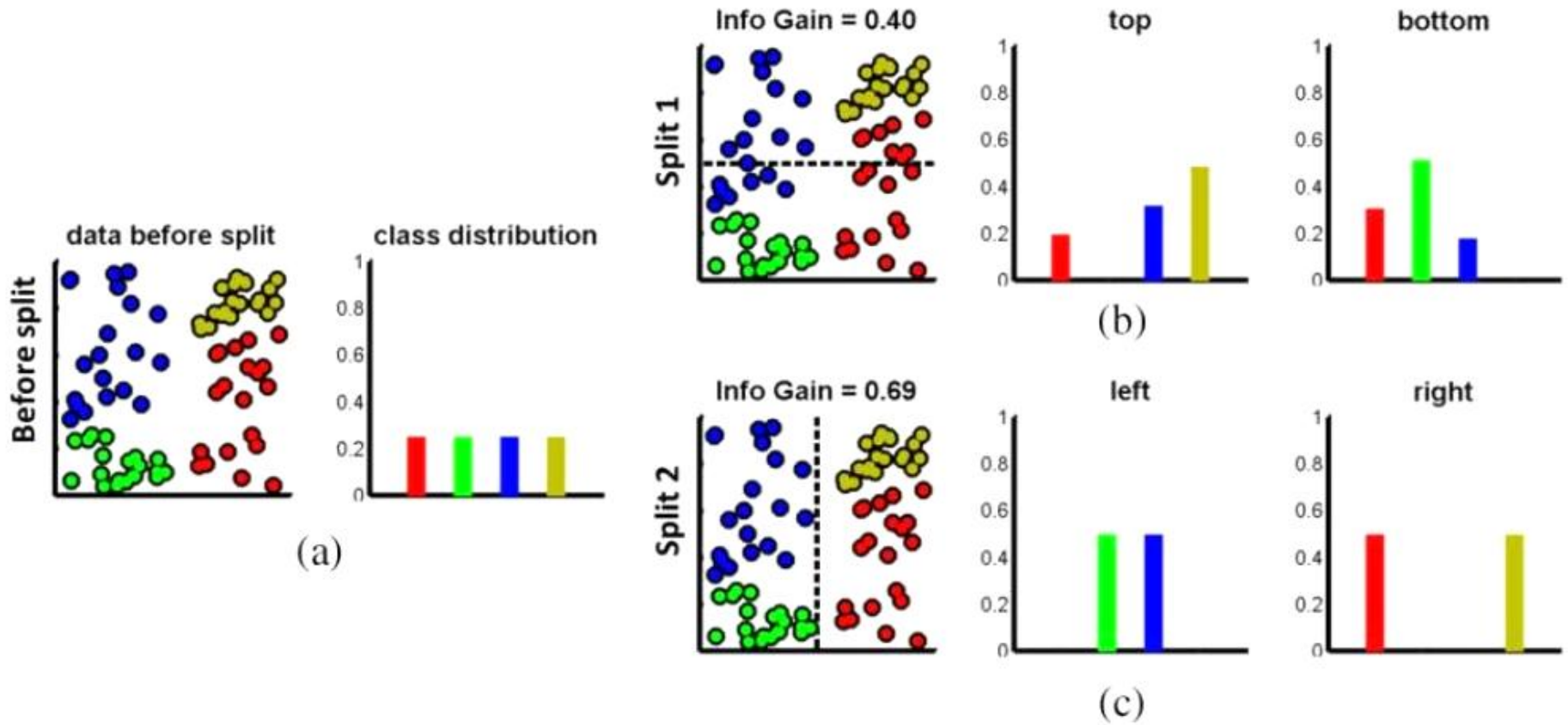
miara (nie)czystości:

- information gain,
- entropy
- Gini Index
- χ^2

- atrybut ma być użyteczny w klasyfikacji
- potrzebna jest miara jakości (im)purity function
- mierzymy stopień separacji obiektów względem klas
- wybieramy atrybut dający „najczystszy” klasowo podział i najmniejsze drzewo



wybór predyktora ⁽⁴⁾



Kryteria oceny podziału

Indeks Gini (algorytmy CART, SPRINT)

Wybieramy atrybut, który minimalizuje indeks Gini

gdzie:

$$\text{gini}(S) = 1 - \sum p_j^2$$

- S – zbiór przykładów należących do n klas
- p_j – względna częstość występowania klasy j w S

Zysk informacyjny (algorytmy ID3, C4.5)

Wybieramy atrybut, który maksymalizuje redukcję entropii

Entropia jest miarą stopnia nieuporządkowania. Im mniejsza wartość entropii, tym większa „czystość” podziału zbioru S na partycje

$$E(A_1, A_2, \dots, A_v) = \sum_{j=1}^v \frac{(s_{1j} + s_{2j} + \dots + s_{mj})}{s} I(s_{1j}, s_{2j}, \dots, s_{mj})$$

Indeks korelacji χ^2 (algorytm CHAID)

Mierzmy korelację pomiędzy każdym atrybutem i każdą klasą (wartością atrybutu decyzyjnego)

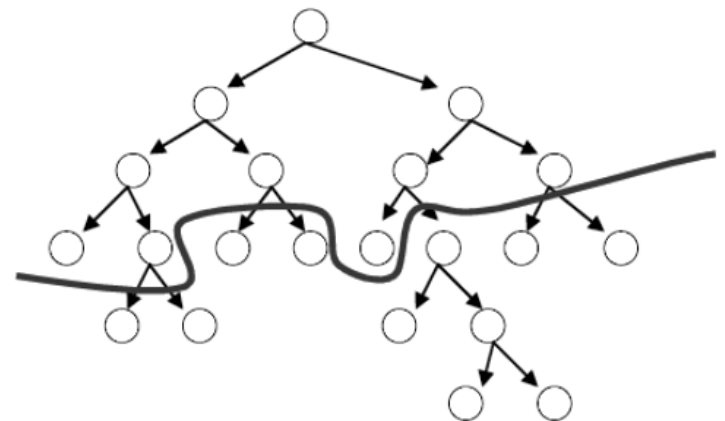
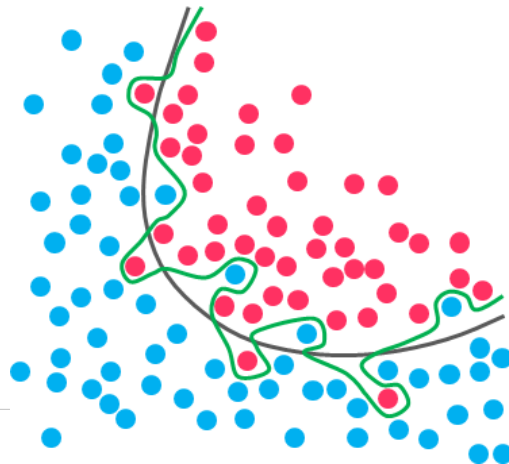
Wybieramy atrybut o maksymalnej korelacji

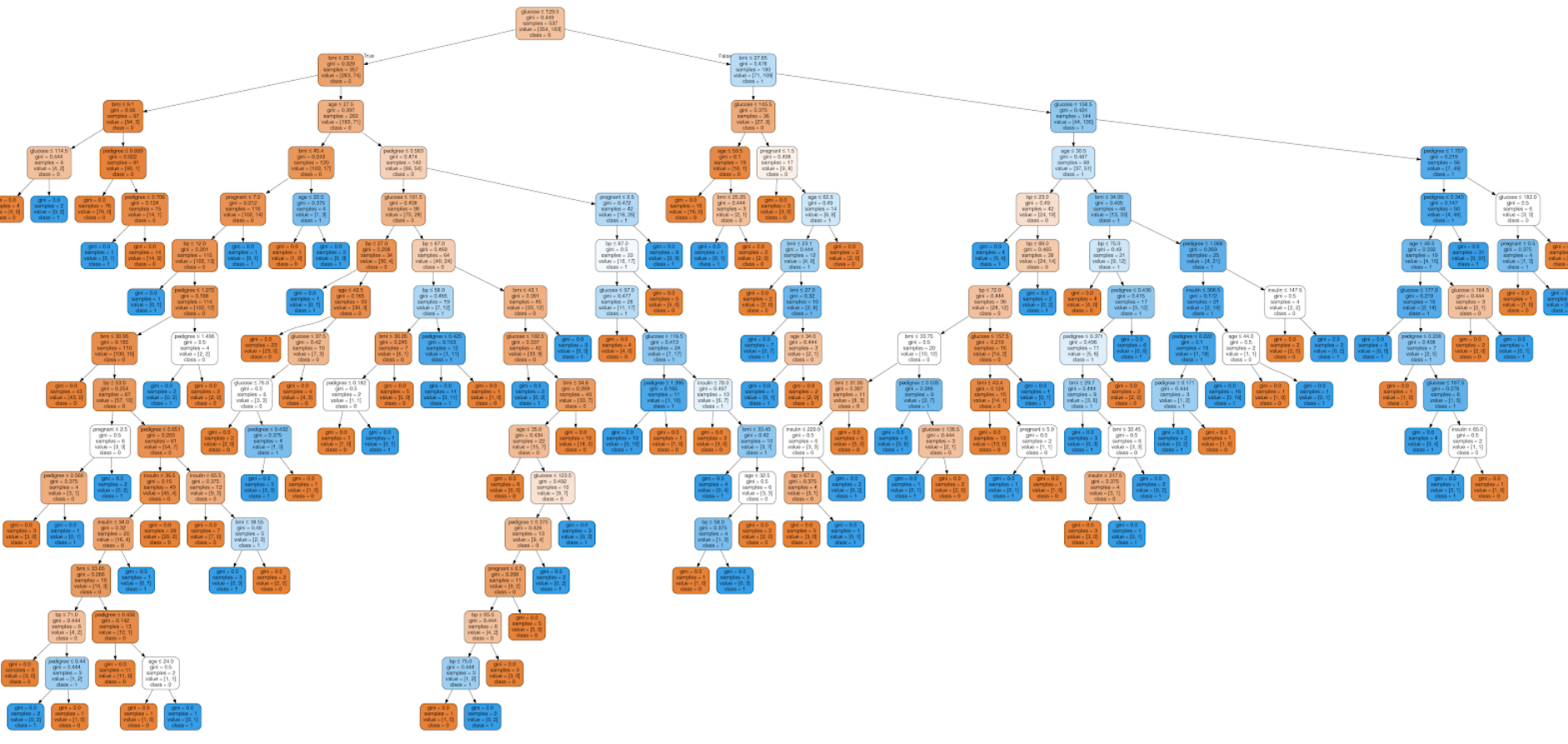
Algorytmy indukcji drzew dążą do możliwie najlepszej klasyfikacji, co prowadzi do przeuczenia (szum i punkty osobliwe)

Przycinanie drzew decyzyjnych - usuwanie mało wiarygodnych gałęzi

- » poprawia efektywność klasyfikacji
- » poprawia zdolność klasyfikatora do klasyfikacji nowych przypadków

Metody przycinania drzew decyzyjnych - bazują najczęściej na miarach statystycznych





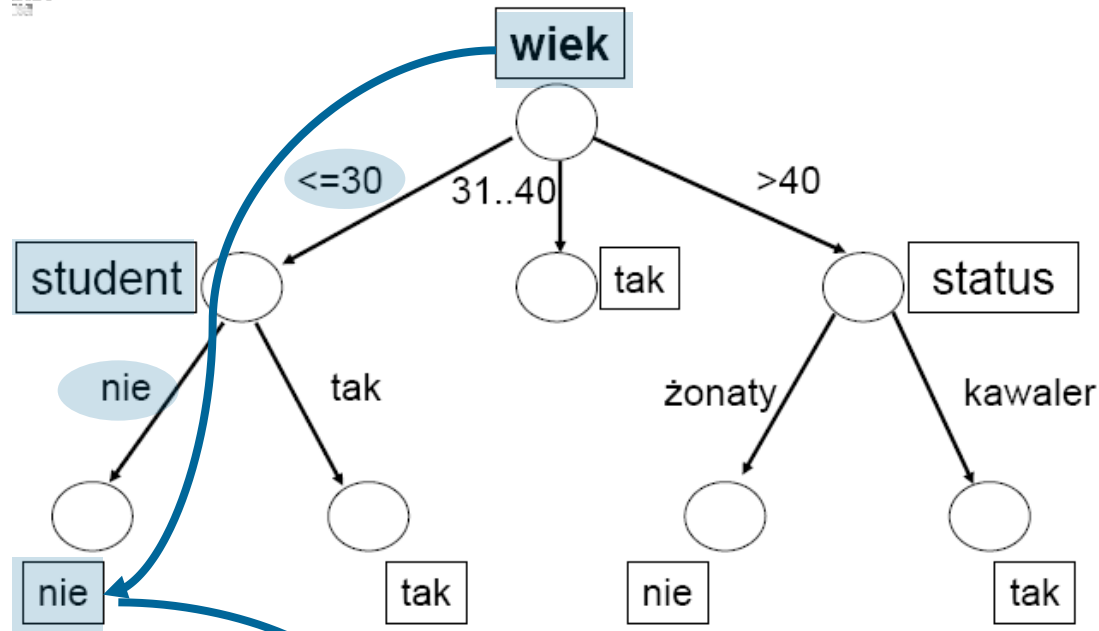
Przykład drzewa bez przycięcia (DataCamp)

Dwa podejścia do problemu przycinania drzew decyzyjnych:

- wstępne przycinanie drzewa decyzyjnego (*prepruning*)
 - » polega na przycięciu drzewa przez wcześniejsze **zatrzymanie** procedury konstrukcji drzewa. Wprowadzamy **warunek stopu**, który wstrzymuje dalsze dzielenie zbioru treningowego na partycje. Przykładowym warunkiem stopu jest przyjęcie **minimalnej liczby elementów** należących do partycji, która podlega dzieleniu.
- przycinanie drzewa po zakończeniu konstrukcji drzewa (*postpruning*)
 - » bazuje na miarach statystycznych

- drzewo decyzyjne można przedstawić w postaci zbioru tzw. **reguł klasyfikacyjnych** postaci *IF-THEN*
- dla każdej **ścieżki** drzewa decyzyjnego, łączącej korzeń drzewa z liściem drzewa tworzymy regułę klasyfikacyjną
- każda **gałąź** tworzy **poprzednik (przesłanki)** reguły klasyfikacyjnej: koniunkcja par $\langle \textit{atrybut}, \textit{wartość} \rangle$
- każdy **liść** tworzy **następnik (konkluzję)** reguły

Ekstrakcja reguł klasyfikacyjnych z drzew decyzyjnych (2)



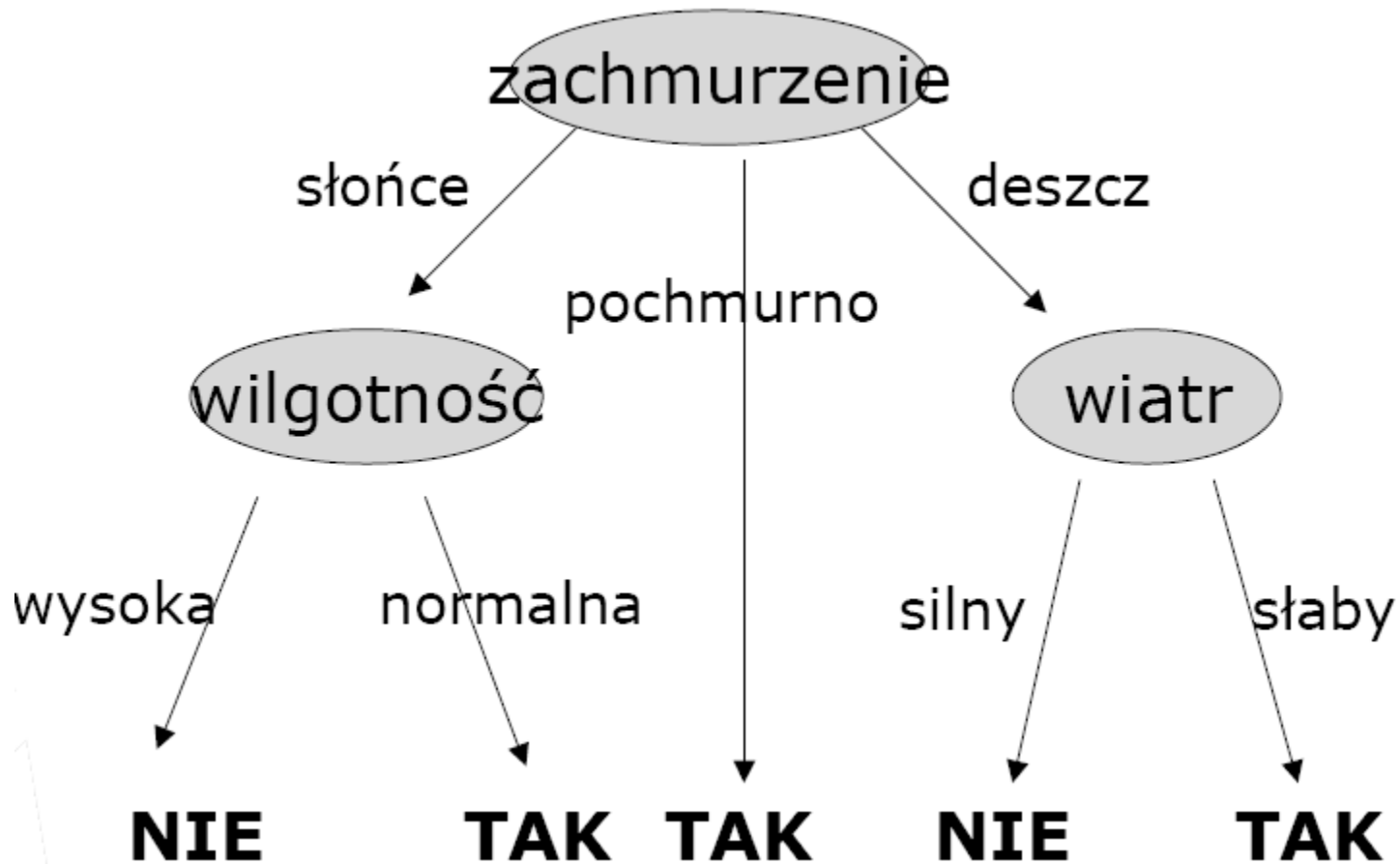
Drzewo decyzyjne można przedstawić w postaci następującego zbioru reguł klasyfikacyjnych:

Reguły:

```

IF wiek='<=30' AND student='nie' THEN kupi_komputer='nie'
IF wiek = '<=30' AND student='tak' THEN kupi_komputer='tak'
IF wiek = '31..40' THEN kupi_komputer = 'tak'
IF wiek = '>40' AND status='żonaty' THEN kupi_komputer = 'nie'
IF wiek = '>40' AND status = 'kawaler' THEN kupi_komputer = 'tak'
    
```

Klasyfikator



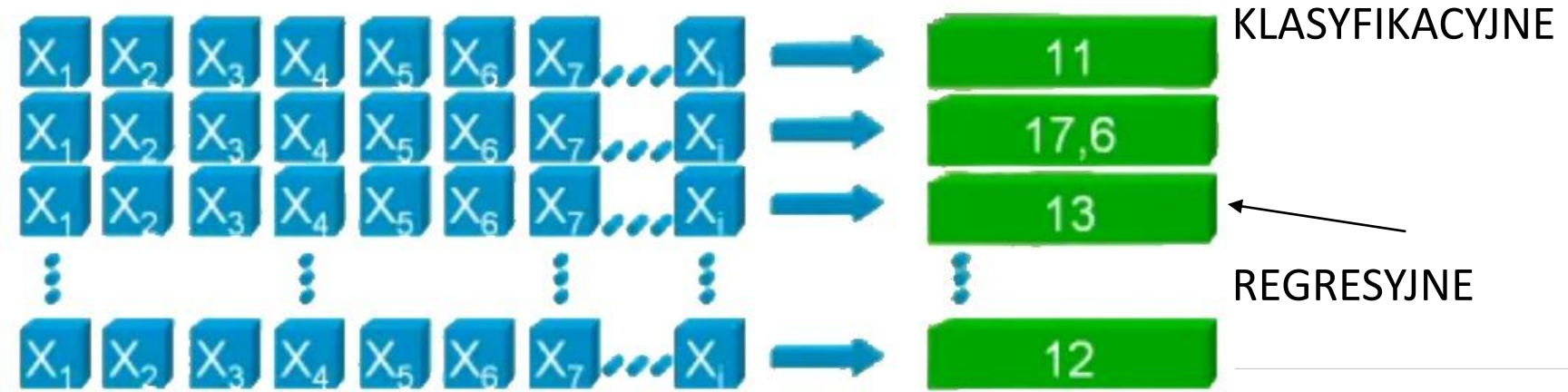
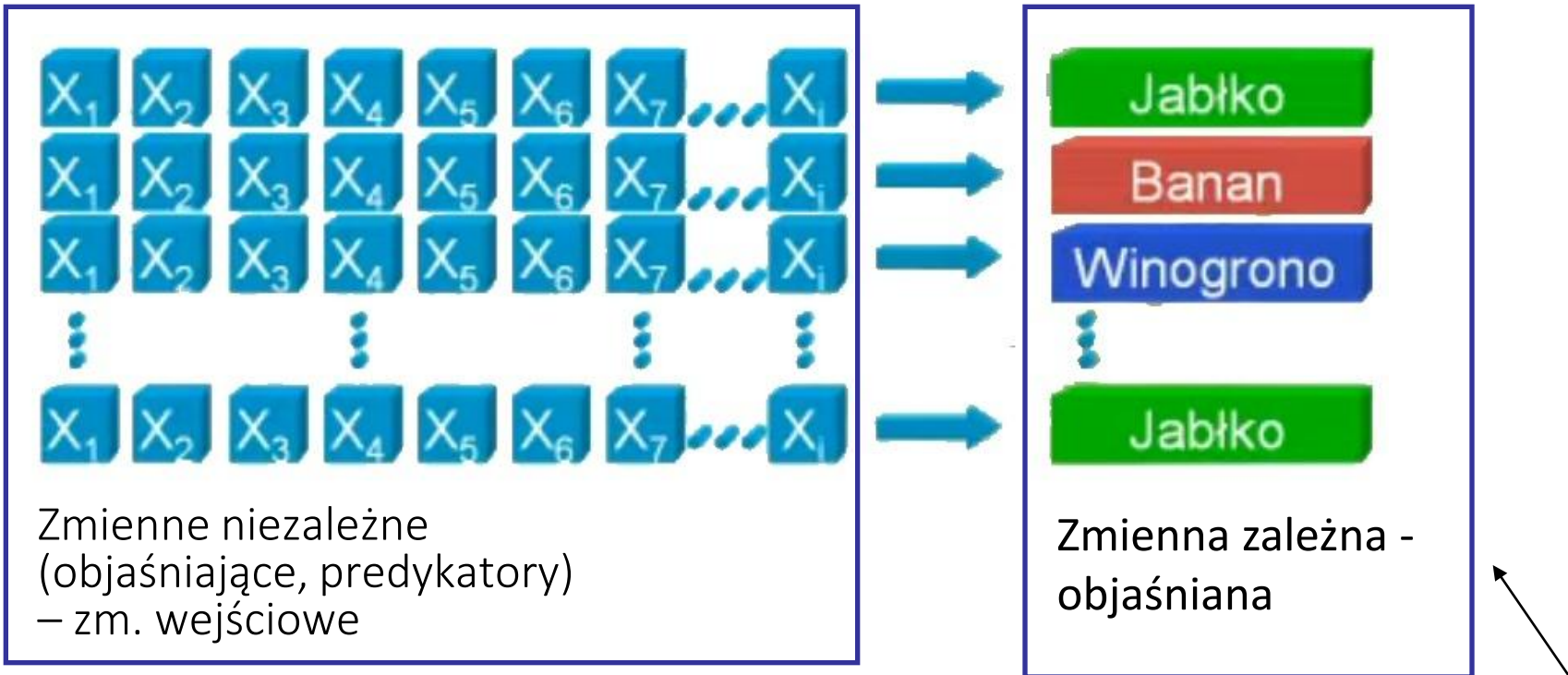
Drzewa – pojęcia podstawowe

- Jeśli zmienna zależna jest wyrażona na skalach słabych (**jakościowych**) to mówimy o drzewach **klasyfikacyjnych**,
- jeśli na mocnych (**ilościowych**), to o drzewach **regresyjnych**.
- Skala zmiennych objaśniających nie ma znaczenia.

Węzły i liście

- Drzewem **binarnym** jest drzewo, w którym z każdego węzła wychodzą dwie krawędzie
- **Liściem** (węzłem końcowym) nazywamy węzeł, z którego nie wychodzą żadne krawędzie
- **Wielkość** drzewa to liczba liści, a **głębokość** drzewa to długość najdłuższej drogi między korzeniem a liściem (liczba krawędzi między tymi dwoma węzłami)

Zmienne w analizie

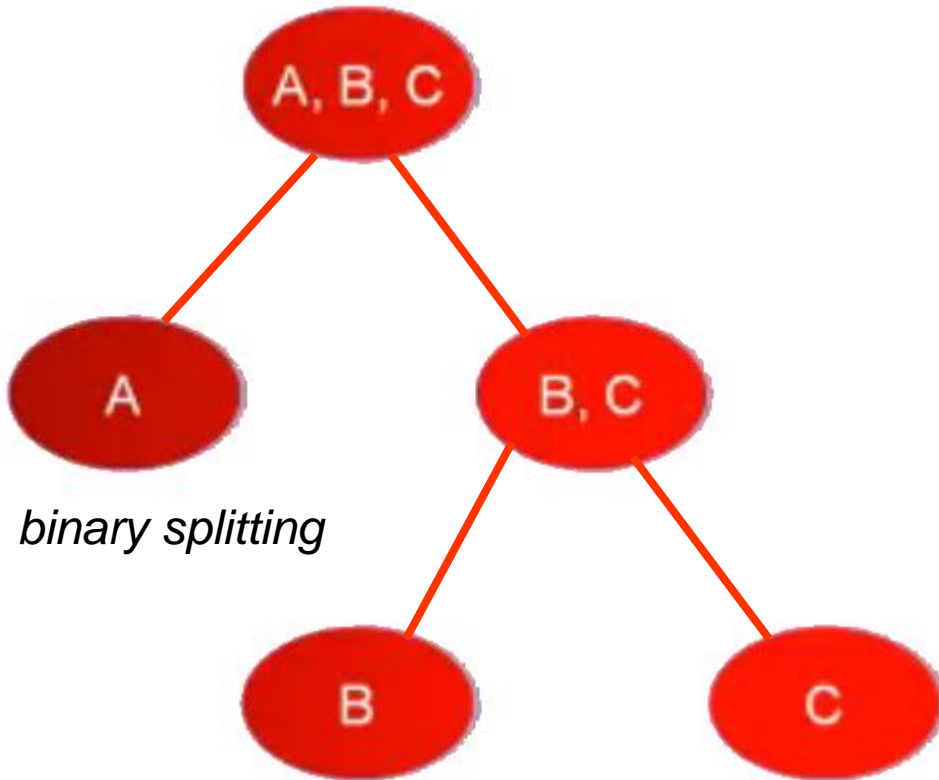


CART (C&RT) i CHAID

Classification & Regression Trees

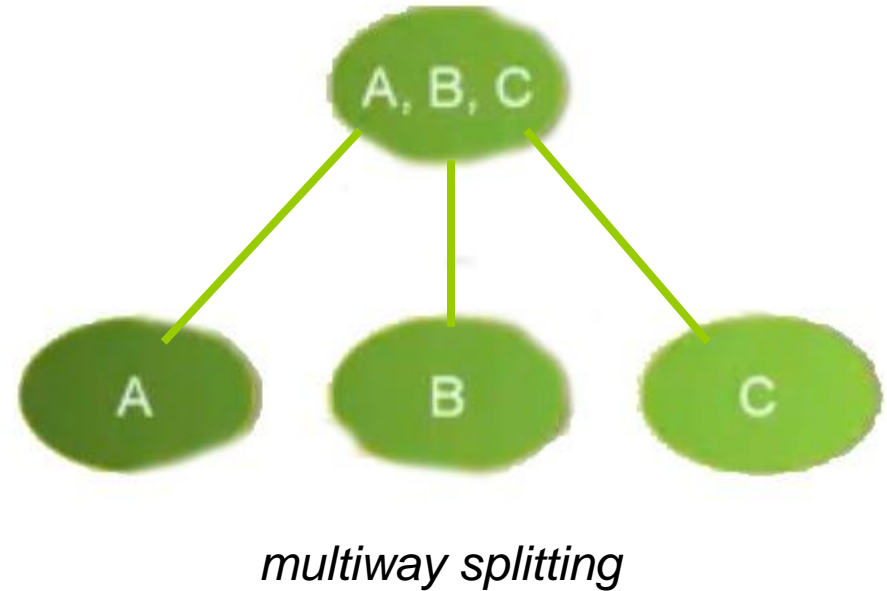
C&RT

CART

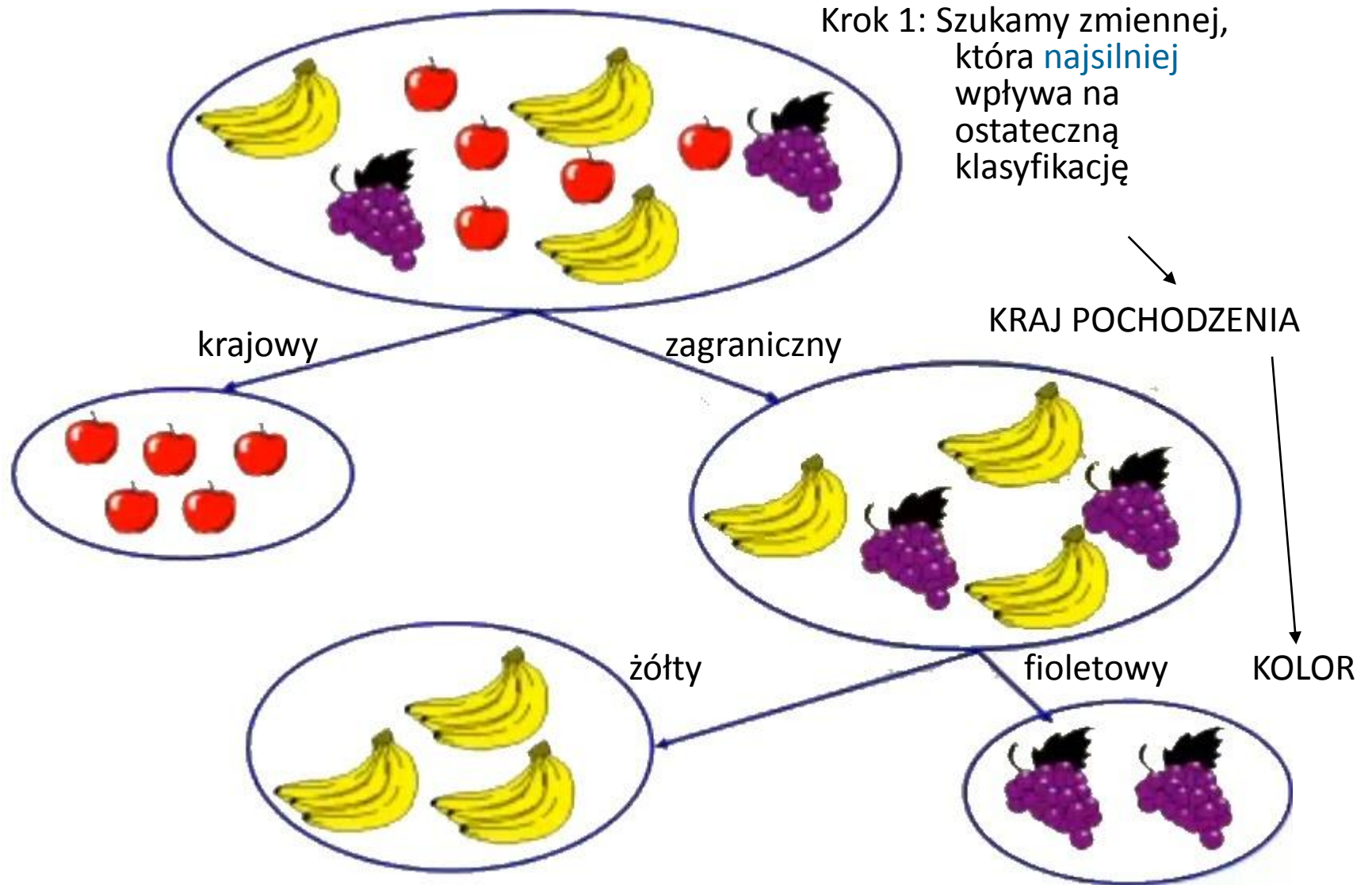


Chi-square Automatic
Interaction Detection

CHAID



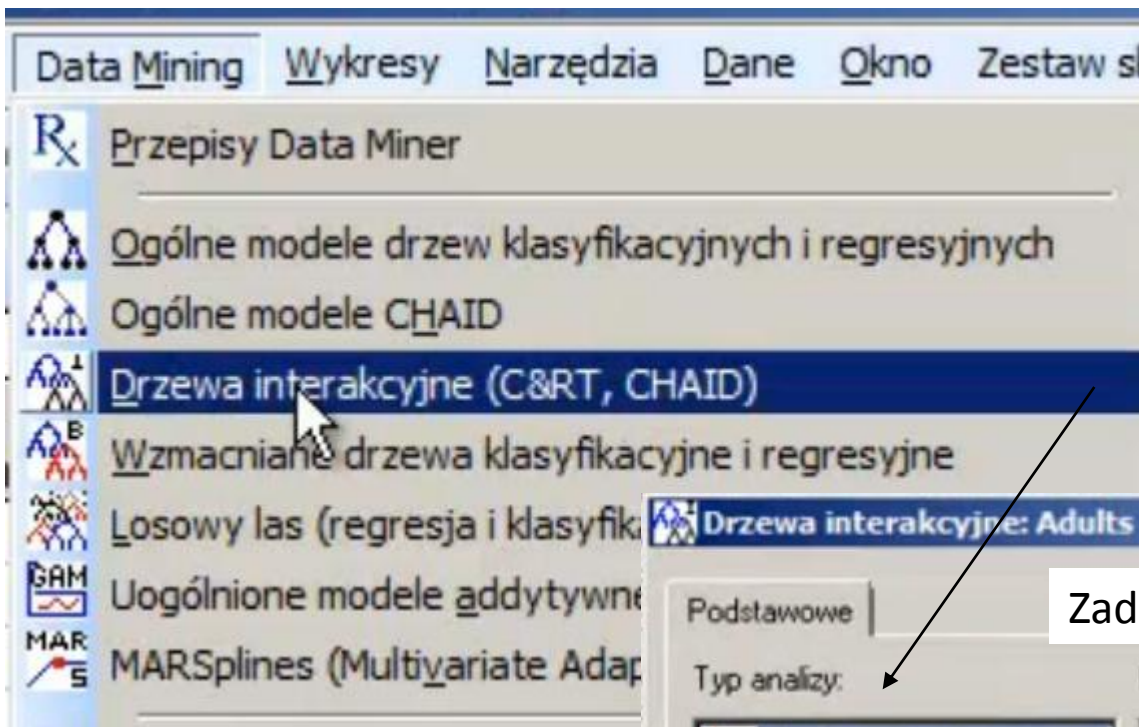
Drzewa klasyfikacyjne



- Wykorzystuje **każdą kombinację** zmiennych ciągłych i kategoryalnych (jakościowych)
- Wybiera **najlepszy podział**
- W kolejnych podziałach mogą być wykorzystywane te same zmienne
- Metoda niewrażliwa na obserwacje odstające
- Obsługuje zbiory obserwacji o złożonej strukturze

Przykład: segmentacja pod kątem dochodów

- Cel: segmentacja ze względu na dochód, charakterystyka segmentów
- Dane zawierają 32 tys przypadków
- Każdy przypadek reprezentuje jedną osobę
- Każda osoba opisana jest przez 11 cech demograficznych oraz zmienną dochód



Drzewa interakcyjne



Zadanie klasyfikacyjne

Opcje rozszerzone interakcyjnego C&RT: Adults

Podstawowe | Klasyfikacja | Stop | Walidacja | Więcej

Zmienne

Zmienna zależna: brak
 Liczności: brak
 Predyktory jakościowe: brak
 Predyktory ilościowe: brak

Kody zm. zależnej: brak

Kody predyktorów:

OK
 Anuluj
 Opcje
 SELECT CASES
 Automatyczna aktualizacja wyników

Wybierz zmienną zależną oraz predyktory jakościowe i ilościowe:

2 - Grupa zawodowa 3 - Wykształcenie 5 - Stan cywilny 6 - Zawód 7 - Związek 8 - Rasa 9 - Płeć 10 - Kraj pochodzenia 11 - Dochód	2 - Grupa zawodowa 3 - Wykształcenie 5 - Stan cywilny 6 - Zawód 7 - Związek 8 - Rasa 9 - Płeć 10 - Kraj pochodzenia 11 - Dochód	1 - Wiek 4 - Liczba lat kształcenia	1 - Wiek 4 - Liczba lat kształcenia
---	---	--	--

Rozwiń Przybliż Rozwiń Przybliż Rozwiń Przybliż Rozwiń Przybliż

Zależna: 11
 Predyktory jakościowe:
 Predyktory ilościowe: 1 4
 Liczności:

Pokazuj tylko zmienne o odpowiedniej skali

OK
 Anuluj
 [Zestawy]...

Włącz opcję "Pokazuj tylko zmienne o odpowiedniej skali" aby na listach, w zależności od potrzeby, pojawiały się tylko zmienne jakościowe albo ilościowe. Naciśnij F1 aby uzyskać więcej informacji.

Wyniki drzew interakc. C&RT: Adults

Menedżer | Podstawowe | Klasyfikacja | Predykcja | Raport

Drzewo (budowa, przycinanie):

- Buduj drzewo
- Buduj i przycinaj drzewo
- Buduj 1 poziom

Przegląd drzewa:

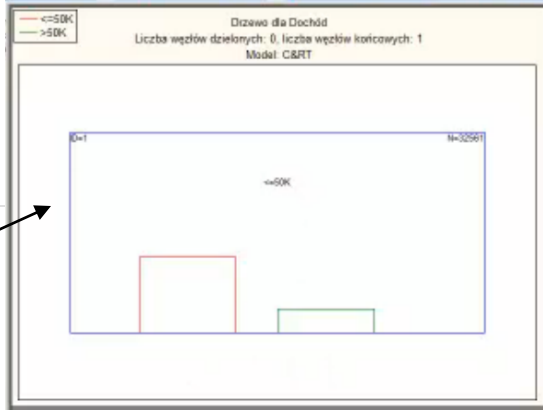
- Przeglądarka
- Przewijalne**
- Drzewo
- Układ drzewa

Węzły i gałęzie:

Węzeł: 1

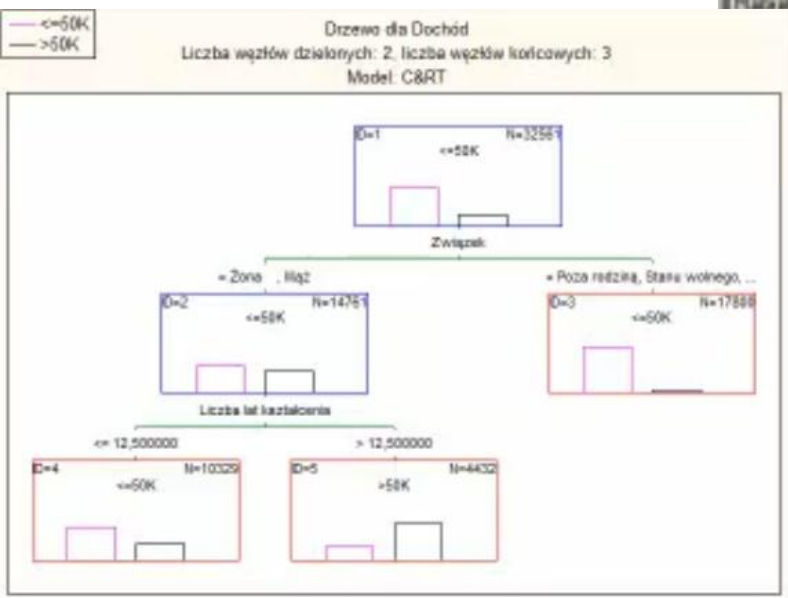
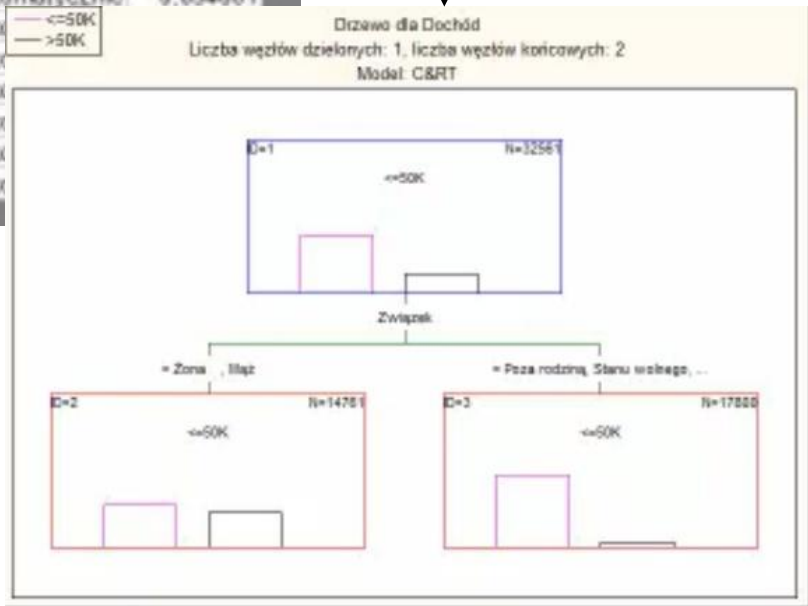
- Buduj gałąź
- Stel. predyktorów
- Warunek podziału
- Dane
- Wybierz zastępcę
- Stat. zastępcy

Zapisz drzewo | Nowe | Zamknij | Opcje



Predykcja dla węzła 1 (Adults)
Model C&RT

Typ podziału	Pograwa Statystyka
Związek	Automatycznie 0,073548
Stan cywilny	Automatycznie 0,069170
Liczba lat kształcenia	Automatycznie 0,039138
Zawód	Automatycznie 0,034661
Wykształcenie	Auto
Wiek	Auto
Płeć	Auto
Grupa zawodowa	Auto
Rasa	Auto
Wzrost	Auto
Przechodzenia	Auto



STATISTICA - przykład drzewa C&RT (inny dobór zmiennych)

Dane: adult (15 zm., * 32561 prz.)

	2	3	4	5	6	7	8	ra
	Work_class	fnlwtg	education	education-num	marital-status	occupation	relationship	ra
1	State-gov	77546	Bachelor	13	Never-married	Adolescent-service-s	Not-in-family	White
2	Self-emp						sband	White
3	Private						t-in-family	White
4	Private						sband	Black
5	Private						fe	Black
6	Private						t-in-family	White
7	Private						sband	Black
8	Self-emp						t-in-family	Black
9	Private						sband	White
10	Private						t-in-family	White
11	Private						sband	Black
12	State-gov						sband	Asia
13	Private						un-ckild	White
14	Private						t-in-family	Black
15	Private						sband	Asian

1 Data Mining Wykresy Narzędzia Dane Okno Pomoc

2 Drzewa interakcyjne: adult

Podstawowe

Typ analizy: Zadanie klasyfikacyjne, Zadanie regresyjne

Metoda budowy modelu: C&RT, CHAID, Wyczerpujący CHAID

Wczytaj drzewo i przejdź do wyników

3 Opcje rozszerzone interakcyjnego C&RT: adult

Podstawowe | Klasyfikacja | Stop | Walidacja | Więcej

Zmienne

Zmienna zależna: bra

Liczności: bra

Predyktory jakościowe: bra

Predyktory ilościowe: bra

Kody zm. zależnej:

Kody predyktorów:

4 Wybierz zmienną zależną oraz predyktory jakościowe i ilościowe:

2 - Work_class	2 - Work_class	1 - Age	1 - Age
4 - education	4 - education	3 - fnlwtg	3 - fnlwtg
6 - marital-status	6 - marital-status	5 - education-num	5 - education-num
7 - occupation	7 - occupation	11 - capital-gain	11 - capital-gain
8 - relationship	8 - relationship	12 - capital-loss	12 - capital-loss
9 - race	9 - race	13 - hours-per-week	13 - hours-per-week
10 - sex	10 - sex		
14 - native-country	14 - native-country		
15 - Income	15 - Income		

Zależna: 15

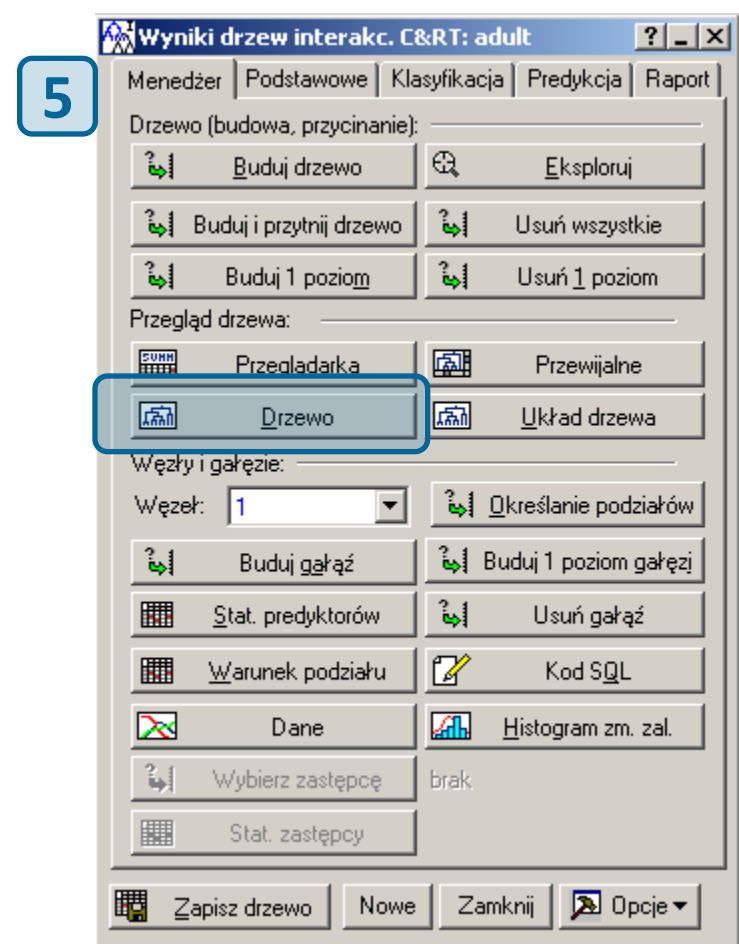
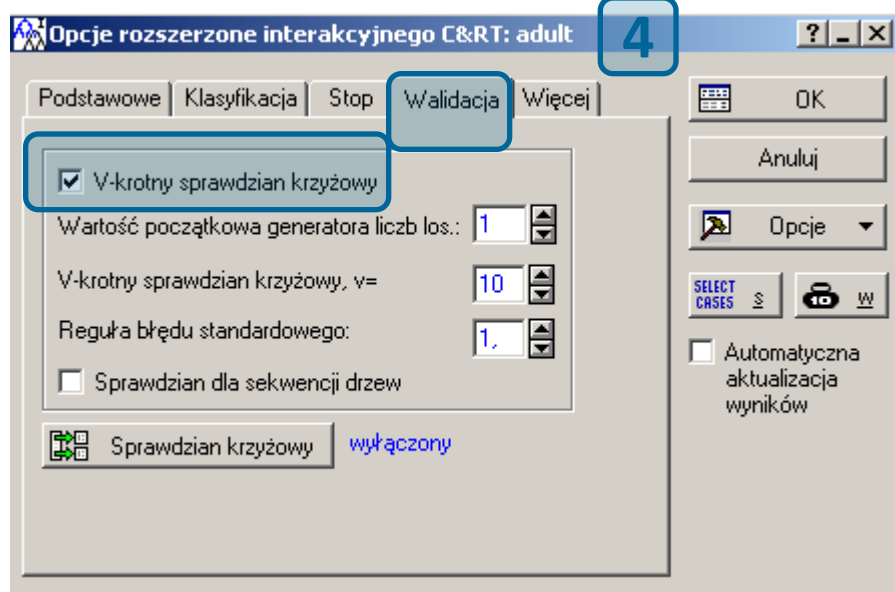
Predyktory jakościowe: 4 6-8

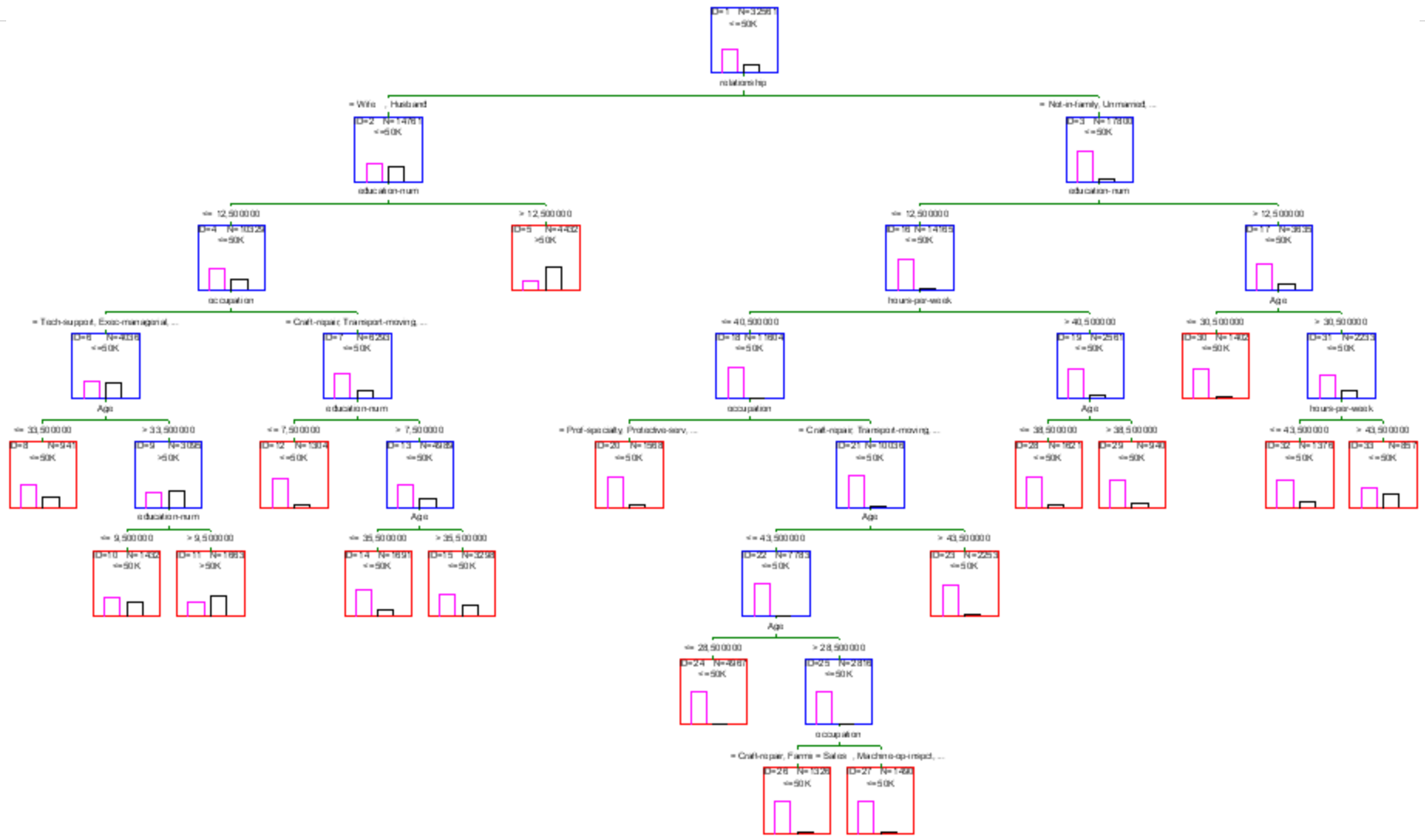
Predyktory ilościowe: 1 5 13

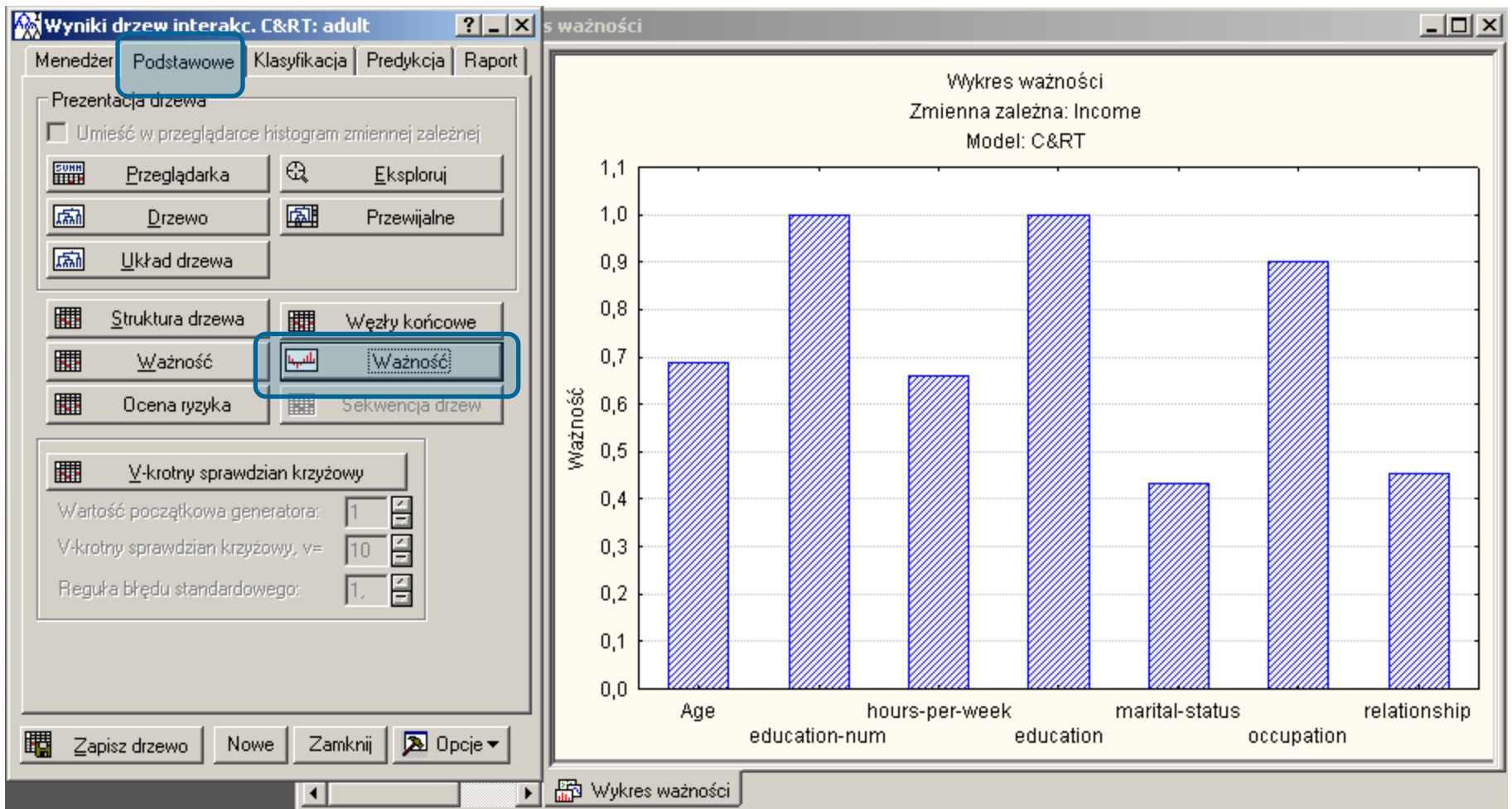
Liczności:

Pokazuj tylko zmienne o odpowiedniej skali

Włącz opcję "Pokazuj tylko zmienne o odpowiedniej skali" aby na listach, w zależności od potrzeby, pojawiały się tylko zmienne jakościowe albo ilościowe. Naciśnij F1 aby uzyskać więcej informacji.







Menedżer | Podstawowe | Klasyfikacja | Predykcja | Raport

Drzewo (budowa, przycinanie):

- Buduj drzewo
- Eksploruj
- Buduj i przycinaj drzewo
- Usuń wszystkie
- Buduj 1 poziom
- Usuń 1 poziom

Przegląd drzewa:

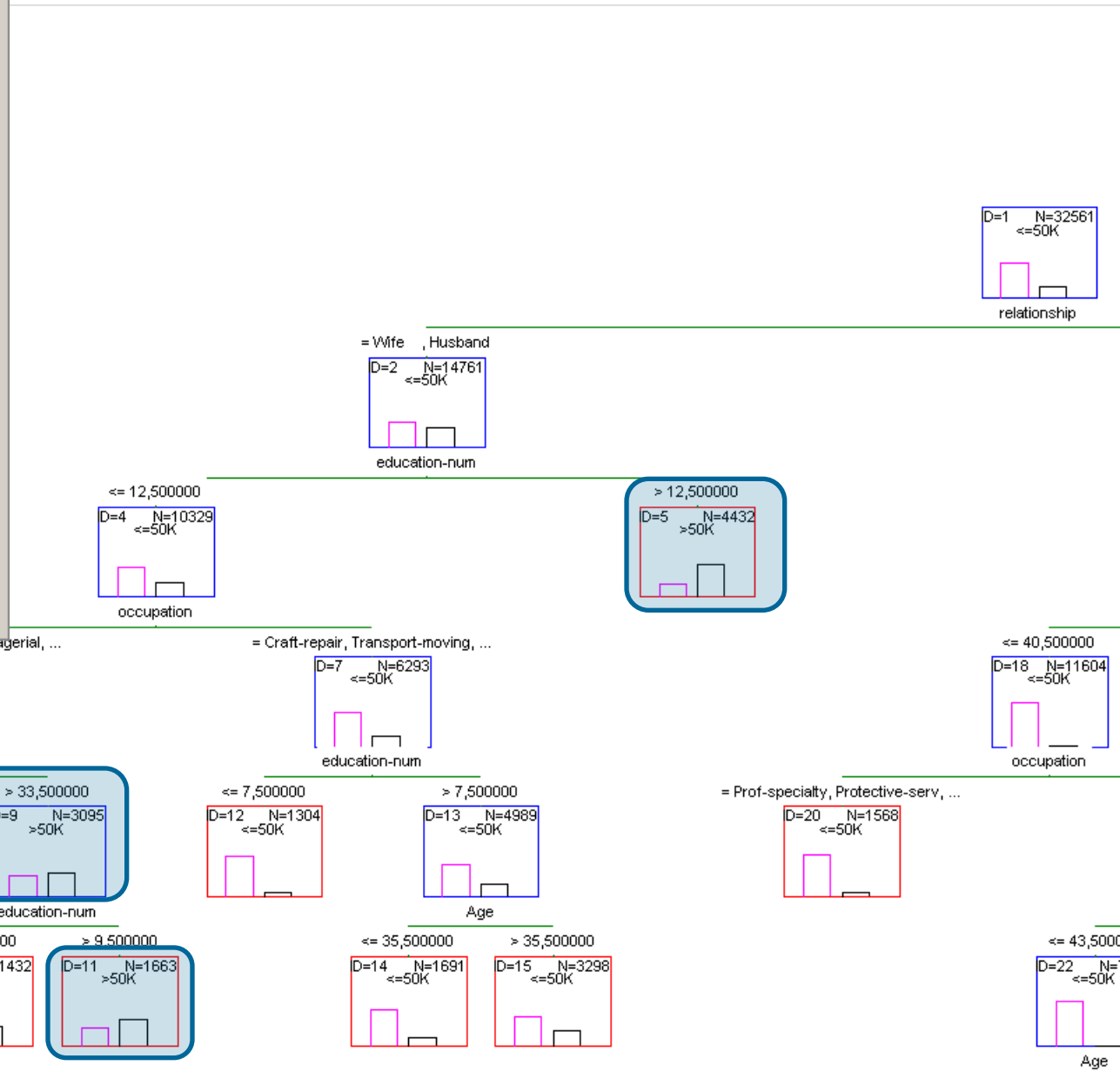
- Przeglądarka
- Przewijalne**
- Drzewo
- Układ drzewa

Węzły i gałęzie:

Węzeł: 1

- Określanie podziałów
- Buduj gałąź
- Buduj 1 poziom gałęzi
- Stat. predyktorów
- Usuń gałąź
- Warunek podziału
- Kod SQL
- Dane
- Histogram zm. zal.
- Wybierz zastępcę
- brak
- Stat. zastępcy

Zapisz drzewo | Nowe | Zamknij | Opcje



Wyniki drzew interakc. C&RT: adult

Menedżer Podstawowe Klasyfikacja Predykcja Raport

Drzewo (budowa, przycinanie):

Buduj drzewo Eksploruj

Buduj i przycinaj drzewo Usuń wszystkie

Buduj 1 poziom Usuń 1 poziom

Przegląd drzewa:

Przeglądarka Przewijalne

Drzewo Układ drzewa

Węzły i gałęzie:

Węzeł: 1 Określanie podziałów

Buduj gałąź Buduj 1 poziom gałęzi

Stat. predyktorów Usuń gałąź

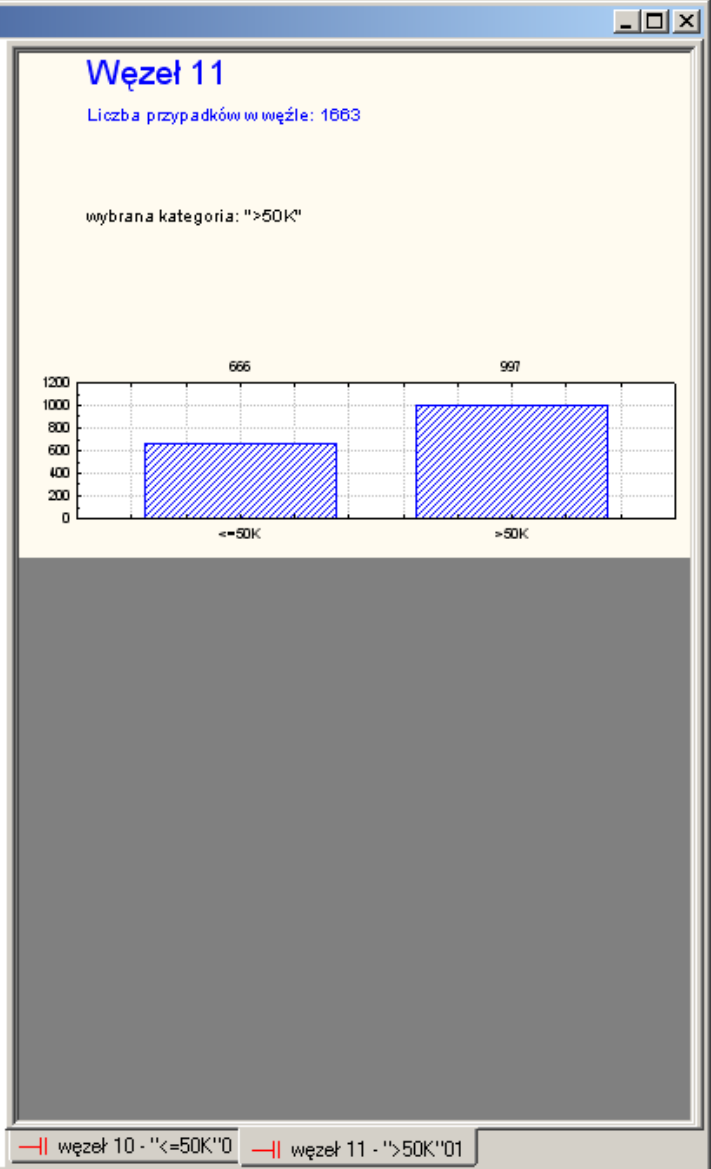
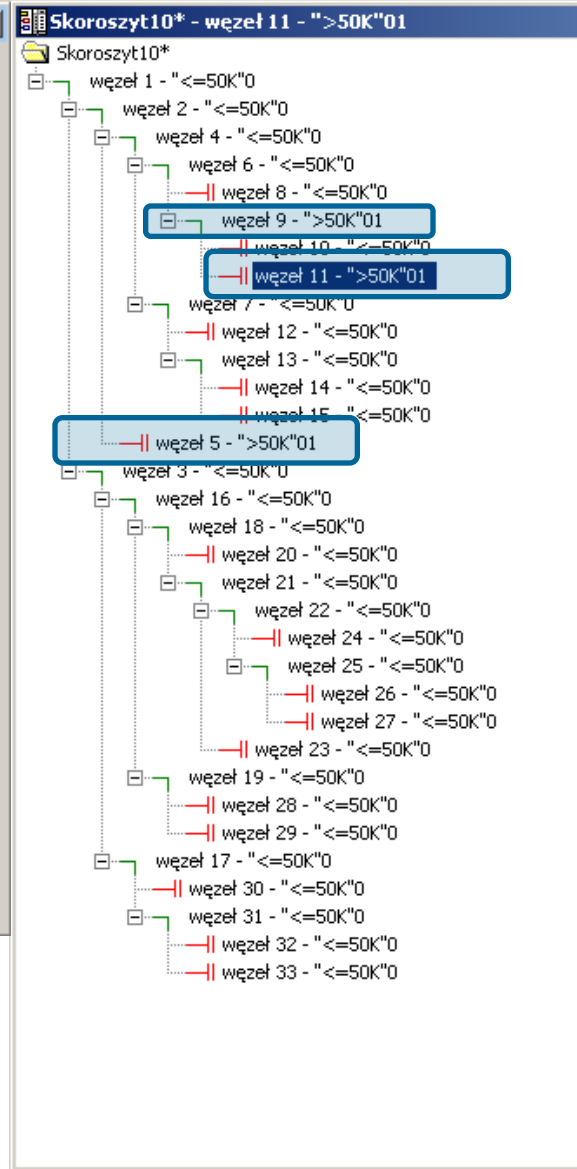
Warunek podziału Kod SQL

Dane Histogram zm. zal.

Wybierz zastępcę brak

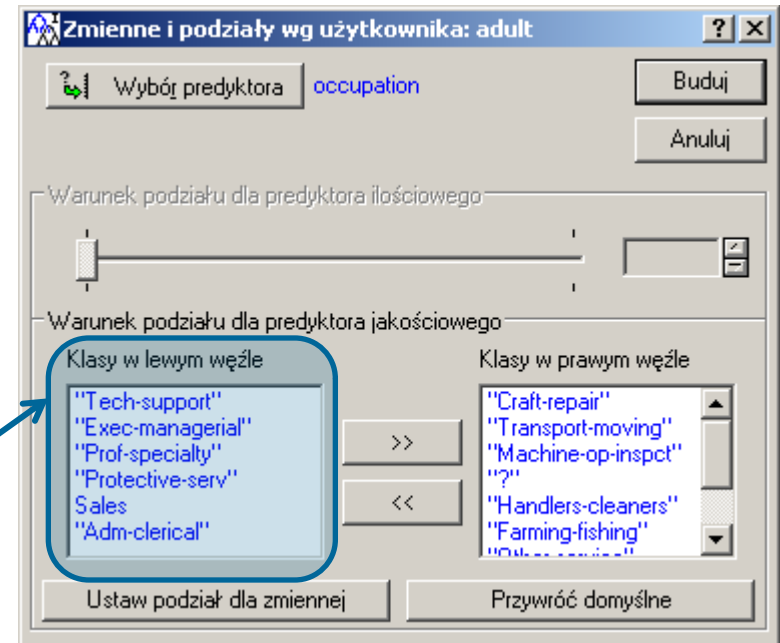
Stat. zastępcy

Zapisz drzewo Nowe Zamknij Opcje



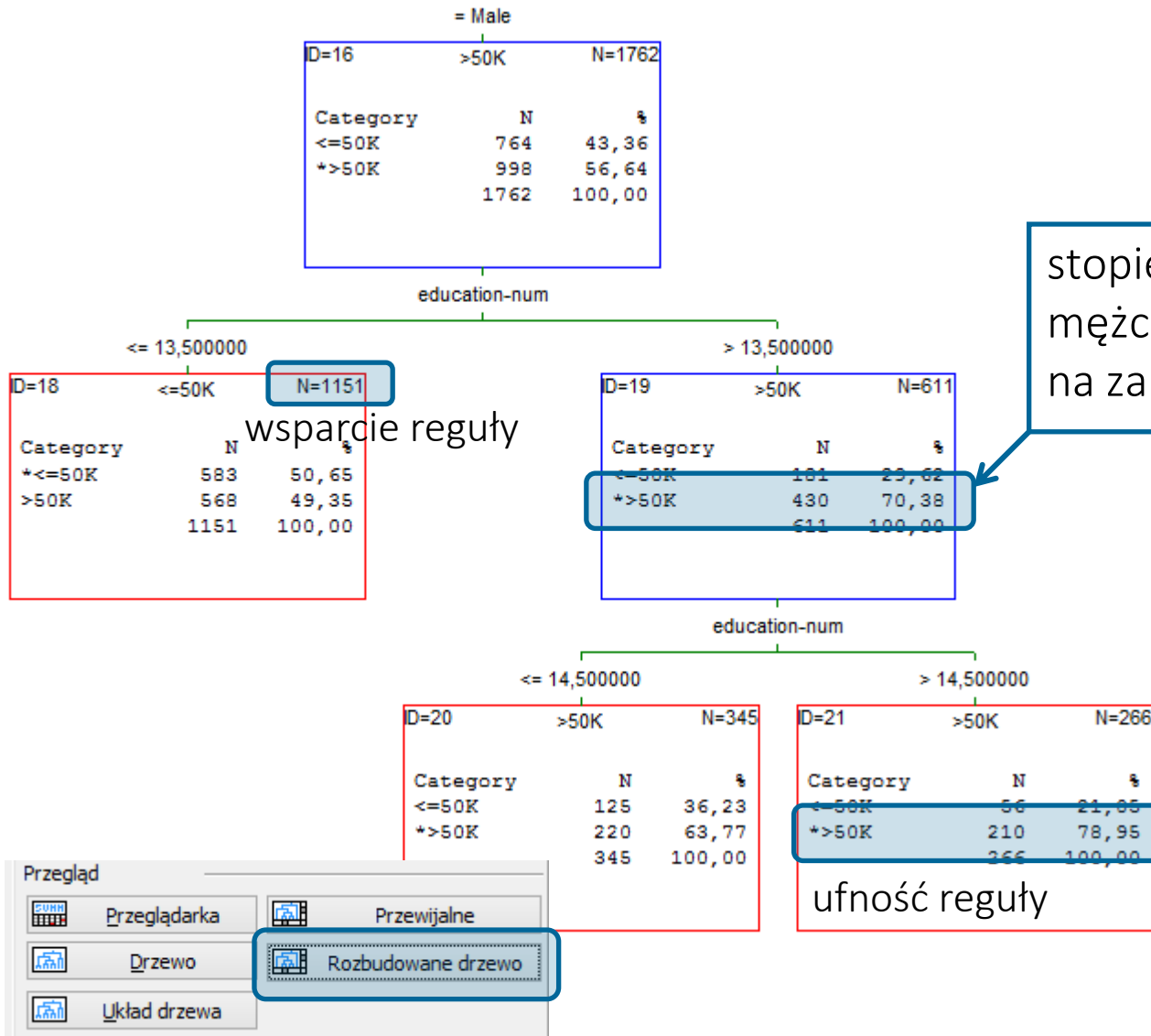
Reguły

- **Jeżeli** osoba pozostaje w związku małżeńskim i jej liczba lat edukacji przekracza 12,5 roku, **wtedy** jej dochód prawdopodobnie przekracza 50 000 \$ (węzeł ID5)
(z prawdopodobieństwem... 72%)
- **Jeżeli** osoba pozostaje w związku małżeńskim, jej liczba lat edukacji nie przekracza 12,5 roku, wykonuje zawód... oraz ma ponad 33,5 lat **wtedy** jej dochód prawdopodobnie przekracza 50 000 \$ (węzeł ID9)
(z prawdopodobieństwem... 53%)
- **Jeżeli** osoba ma ponad 33,5 lat, pozostaje w związku małżeńskim, liczba lat jej edukacji mieści się w przedziale 9,5 do 12,5 lat, wykonuje zawód... **wtedy** jej dochód prawdopodobnie przekracza 50 000 \$ (węzeł ID11)
(z prawdopodobieństwem... 60%)



*przycisk:
„określanie podziałów”*

Pewność reguł (ufność reguły, dokładność)



stopień magistra daje mężczyznom 70,38% szans na zarobki pow.50tys

studia podyplomowe i doktorat podnosi szansę na zarobki pow.50tys o ponad 8 punktów procentowych

wsparcie reguły

ufność reguły

Przegląd

- Przeglądarka
- Przewijalne
- Drzewo
- Rozbudowane drzewo**
- Układ drzewa

Wsparcie i Ufność

jest bardzo mało kobiet wysoko wykształconych

D=1 <=50K N=9950		
Category	N	%
*<=50K	7540	75,78
>50K	2410	24,22
	9950	100,00

Wsparcie =
 $217/9950$
 = 2,2%

większość kobiet, nawet bardzo dobrze wykształconych, nie zarabia pow. 50 tys.

Female

D=17 <=50K N=697		
Category	N	%
*<=50K	518	74,32
>50K	179	25,68
	697	100,00

education-num

<= 13,500000

> 13,500000

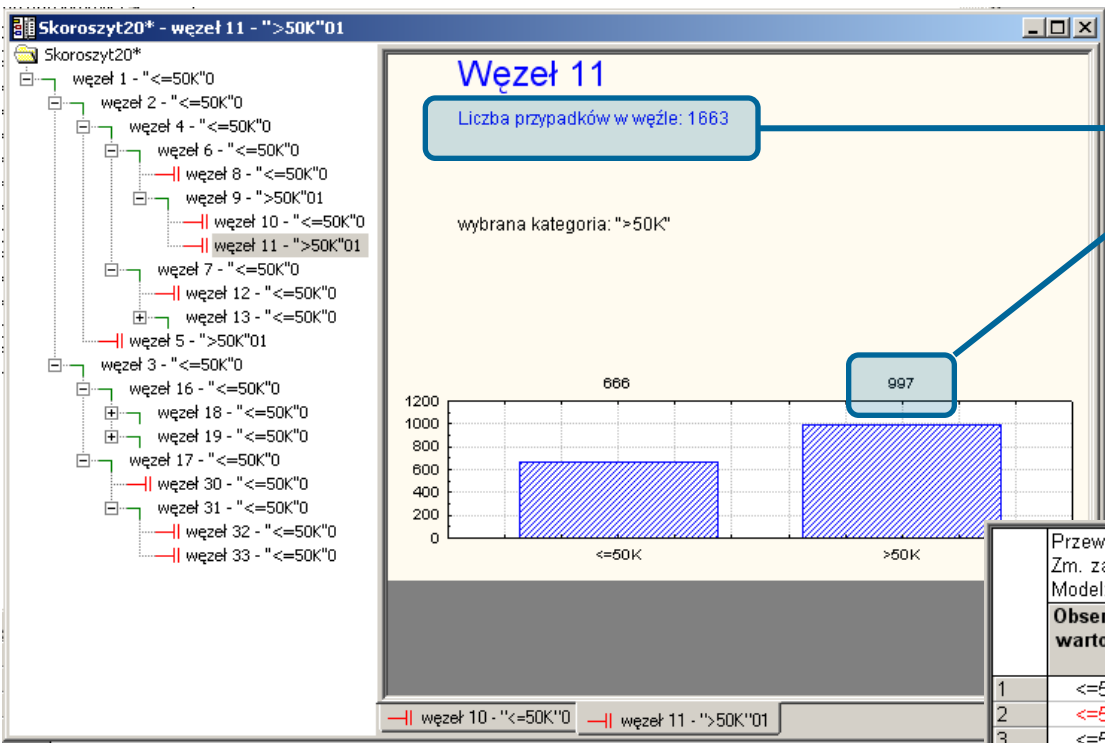
D=22 <=50K N=480		
Category	N	%
*<=50K	390	81,25
>50K	90	18,75
	480	100,00

D=23 <=50K N=217		
Category	N	%
*<=50K	128	59,99
>50K	89	41,01
	217	100,00

Ufność =
 41,01%

kobiety ze stopniem magistra (i wyżej) mają jedynie 41% szans na zarobki >50K

Wsparcie (pokrycie) reguły / ufność (dokładność)



$$ID11: 997/1633 = 0,5995$$

ufność reguły
 $N_{konkluzji} / N_{węzła}$

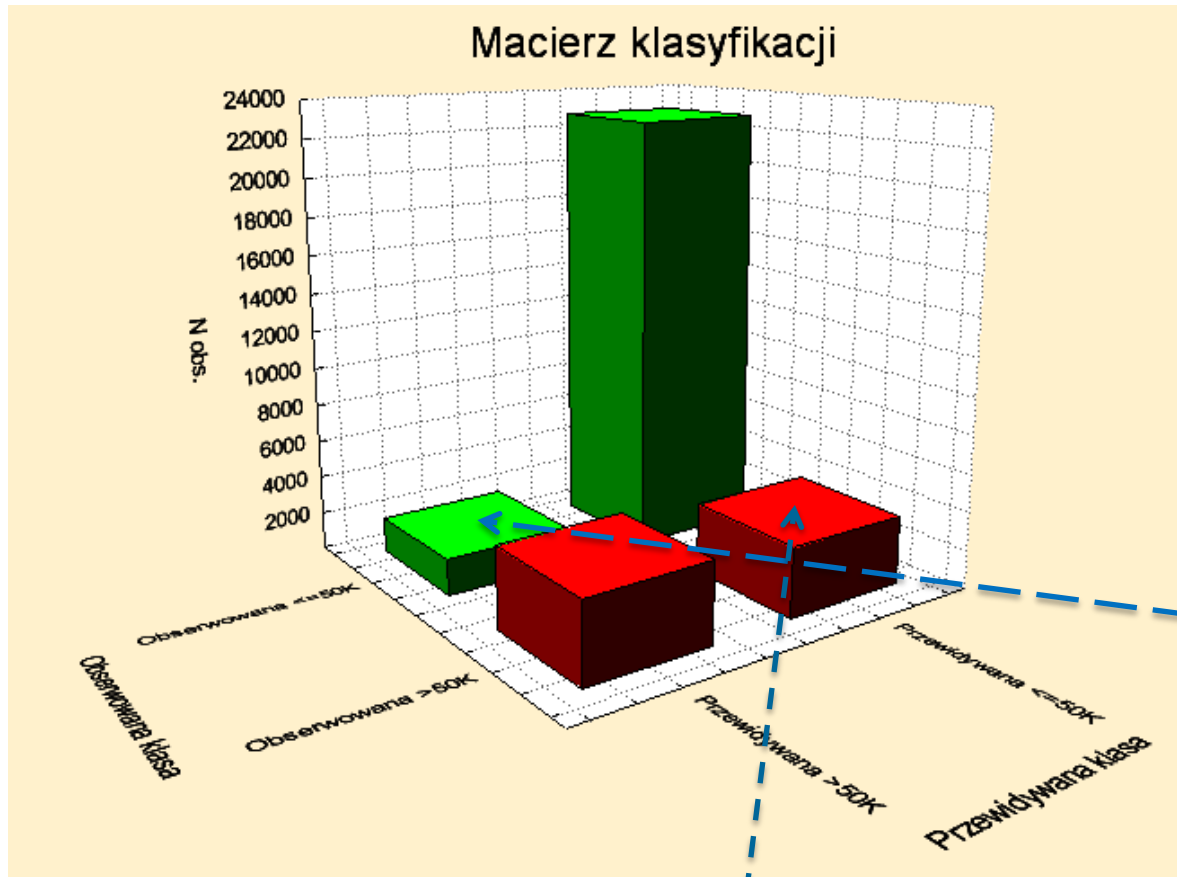
wsparcie reguły:

$$N_{węzła} / N_{zbioru}$$

$$ID11: 1633/32561=5\%$$

	Przewidywane (adult)	Zm. zal.:Income	Model: C&RT		
	Obserw. wartość	Przewid. wartość	Prawdopodobieństwo <=50K	Prawdopodobieństwo >50K	Końcowe węzły
1	<=50K	<=50K	0,814680	0,185320	32
2	<=50K	>50K	0,276399	0,723601	5
3	<=50K	<=50K	0,993289	0,006711	27
4	<=50K	<=50K	0,900307	0,099693	12
5	<=50K	>50K	0,276399	0,723601	5
6	<=50K	>50K	0,276399	0,723601	5
7	<=50K	<=50K	0,967155	0,032845	23
8	>50K	<=50K	0,560754	0,439246	10
9	>50K	<=50K	0,577596	0,422404	33
10	>50K	>50K	0,276399	0,723601	5
11	>50K	>50K	0,400481	0,599519	11
12	>50K	>50K	0,276399	0,723601	5
13	<=50K	<=50K	0,939372	0,060628	30
14	<=50K	<=50K	0,942011	0,057989	28
15	>50K	<=50K	0,660703	0,339297	15
16	<=50K	<=50K	0,900307	0,099693	12
17	<=50K	<=50K	0,997584	0,002416	24
18	<=50K	<=50K	0,993289	0,006711	27

Macierz klasyfikacji



Ile razy model się pomylił?
Pojęcie „kosztu”

false positives

Macierz klasyfikacji (adult)
Zm. zal.: Income
Model: C&RT

Obserw.	Przewidywana <=50K	Przewidywana >50K	Łącznie w wierszu
Liczba	<=50K 22829	1891	24720
Procent z kolumny	86.26%	31.03%	
Procent z wiersza	92.35%	FP 7.65%	
Procent z ogółu	70.11%	5.81%	75.92%
Liczba	>50K 3637	4204	7841
Procent z kolumny	13.74%	68.97%	
Procent z wiersza	16.38%	53.62%	
Procent z ogółu	1.17%	12.91%	24.08%
Liczba	Ogół grup 26466	6095	32561
Procent łącznie	81.28%	18.72%	

false negatives

FN

- Algorytm drzew klasyfikacyjnych
- Zmienne ilościowe dzielone są na 10 kategorii, zmienne jakościowe obsługiwane w sposób naturalny
- Wyszukiwanie par kategorii podobnych do siebie ze względu na zmienną zależną
- Test χ^2

- Co wpływa na **skłonność zakupu samochodu nowego bądź używanego?**
- Wybór jednego z 12 profili aut o porównywalnej cenie (połowa z nich używane, połowa – nowe)
- 1200 ankietowanych,
- dane demograficzne + wybór

Wyświetlony pasek menu i paski narzędzi w programie Data Mining. W menu widoczne są opcje: Wyświetl, Wstaw, Format, Statystyka, Data Mining, Wykresy, Narzędzia, Dane, Okno, Zestaw skoringowy, Pomoc. Paski narzędzi zawierają ikony do kopiowania, wklejania, usuwania, drukowania oraz narzędzia do formatowania tekstu i tabel.

	1 samochód	2 model	3 kraj pochodzenia marki	4 niemieckie - pozostałe	5 prawo jazdy	6 auto - badany	7 auto - rodzice	8 płeć	9 miejscowość	10 województwo	11 tryb audycji
1	używany	VW	Niemcy	Niemcy	tak	tak					
2	używany	VW	Niemcy	Niemcy	tak	tak					
3	używany	VW	Niemcy	Ni							
4	używany	Toyota	Japonia	in	tak	tak	tak	kobieta	do 50 tys.	małopolskie	zacczna
5	używany	Audi	Niemcy	Ni	tak	tak	tak	mężczyzna	wieś	śląskie	zacczna
6	używany	Toyota	Japonia	in	tak	tak	tak	mężczyzna	do 50 tys.	śląskie	zacczna
7	używany	Audi	Niemcy	Ni	tak	tak	tak	kobieta	> 200 tys.	lubelskie	dzienne
8	nowy	Fiat	inny kraj	in	tak	tak	tak	mężczyzna	> 200 tys.	podkarpackia	dzienne
9	nowy	Fiat	inny kraj	in	tak	tak	tak	mężczyzna	100-200 tys.	śląskie	dzienne
10	nowy	Fiat	inny kraj	in	tak	tak	tak	mężczyzna	> 200 tys.	małopolskie	dzienne
11	używany	VW	Niemcy	Ni	tak	tak	tak	kobieta	> 200 tys.	małopolskie	dzienne
12	używany	Audi	Niemcy	Ni	tak	tak	tak	kobieta	do 50 tys.	mazowieckie	dzienne

Wyświetlony panel boczny z menu. Wybrana opcja: **Drzewa interakcyjne (C&RT, CHAID)**. Inne widoczne opcje to: Przepisy Data Miner, Ogólne modele drzew klasyfikacyjnych i regresyjnych, Ogólne modele CHAID, Wzmacnianie drzew klasyfikacyjnych i regresyjnych, Losowy las (regresja i klasyfikacja), Uogólnione modele addytywne, MARSplines (Multivariate Adaptive Regression Splines).

Wyświetlony dialogowy okno konfiguracji: **Drzewa interakcyjne: Intencje_zakupowe_ankieta_los**.
 Typ analizy: Zadanie klasyfikacyjne (wybrane), Zadanie regresyjne.
 Metoda budowy modelu: C&RT, CHAID (wybrane), Wyczerpujący CHAID.
 Przyciski: OK, Anuluj, Opcje, Otwórz dane.
 Na dole: Wczytaj drzewo i przejdź do wyników, SELECT CASES, W.

Wybierz zmienną zależną oraz predyktory jako

- samochód	1 - samochód
2 - model	2 - model
3 - kraj pochodzenia marki	3 - kraj pochodzenia marki
4 - niemieckie - pozostałe	4 - niemieckie - pozostałe
5 - prawo jazdy	5 - prawo jazdy
6 - auto - badany	6 - auto - badany
7 - auto - rodzice	7 - auto - rodzice
8 - płeć	8 - płeć
9 - miejscowosc	9 - miejscowosc
10 - województwo	10 - województwo
11 - tryb studiów	11 - tryb studiów

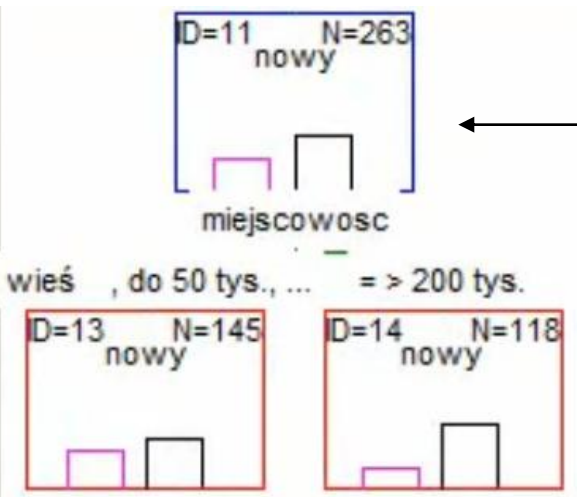
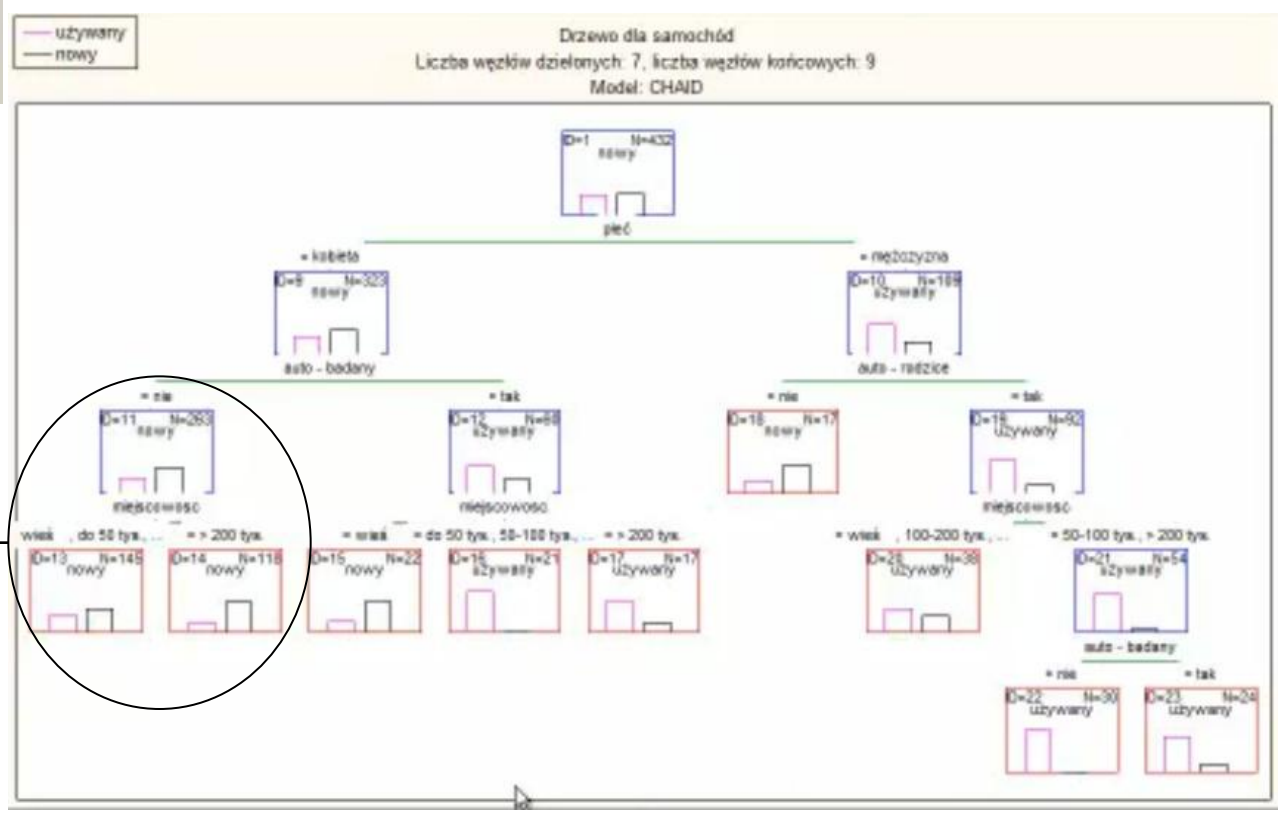
Rozwiń Przybliż Rozwiń Przybliż

Zależna: Predyktory jakościowe:

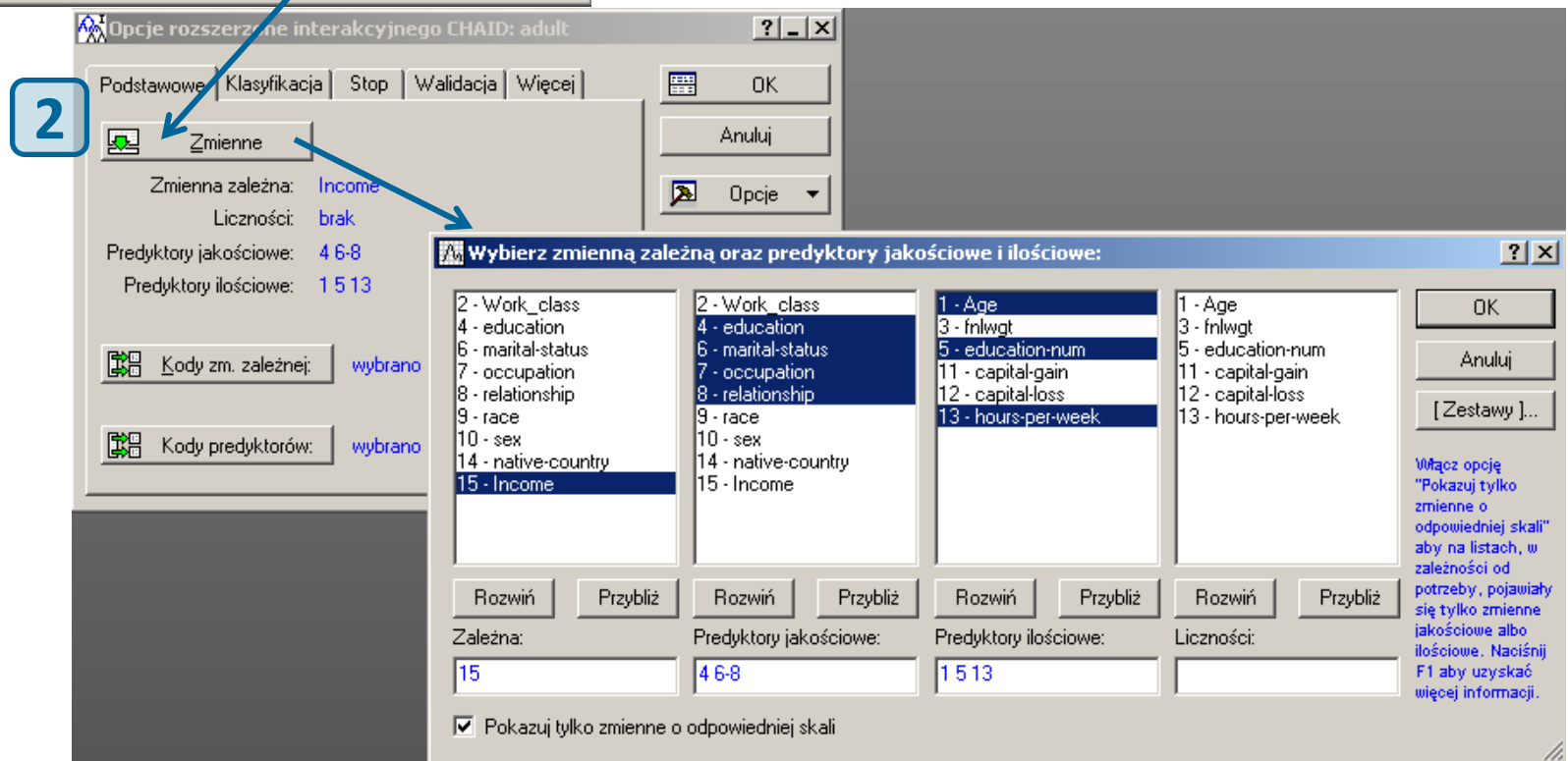
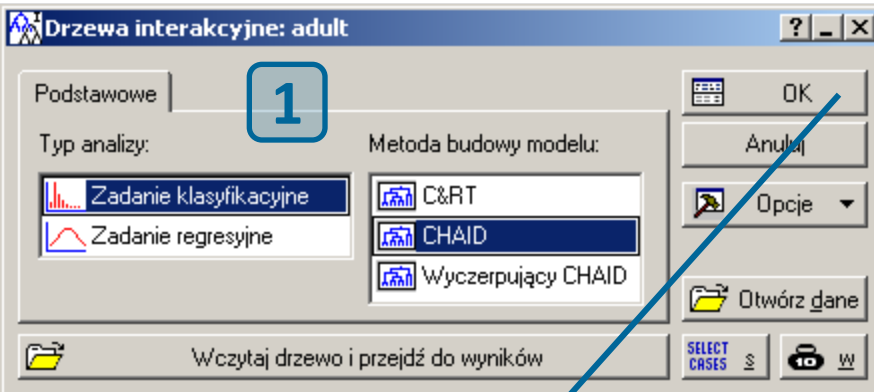
Model: CHAID

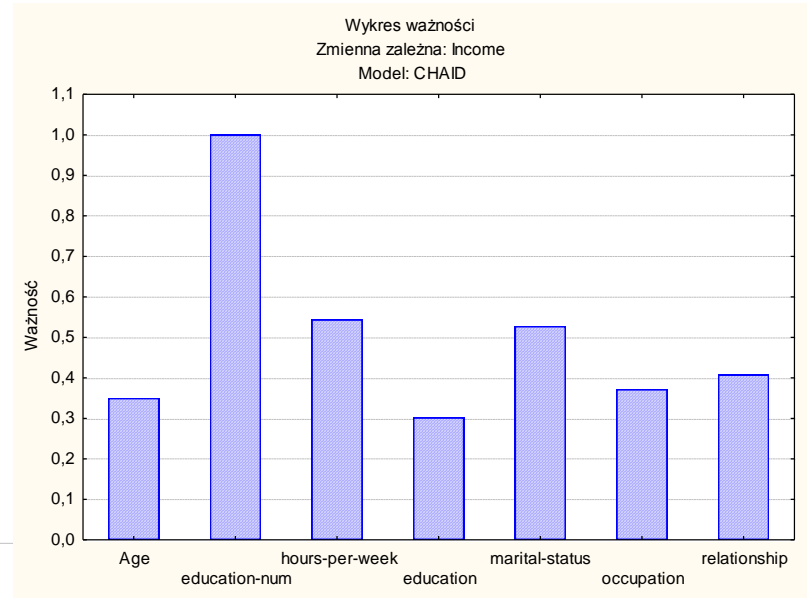
	Liczba węzły	Typ podziału	chi-kwadrat Statystyka	Stopnie Swobody	Skorygowane p
płeć ➕➔	2	Automatycznie	27,13573	1,000000	0,000000
prawo jazdy	2	Automatycznie	21,86719	1,000000	0,000003
auto - badany	2	Automatycznie	17,78424	1,000000	0,000025
tryb studiów	3	Automatycznie	28,19740	2,000000	0,000001
województwo	2	Automatycznie	12,80863	1,000000	0,022082
miejscowosc	2	Automatycznie	11,17302	1,000000	0,013279
auto - rodzice	2	Automatycznie	3,89359	1,000000	0,048471

Zmienną decyzyjną najsilniej różnicuje płeć

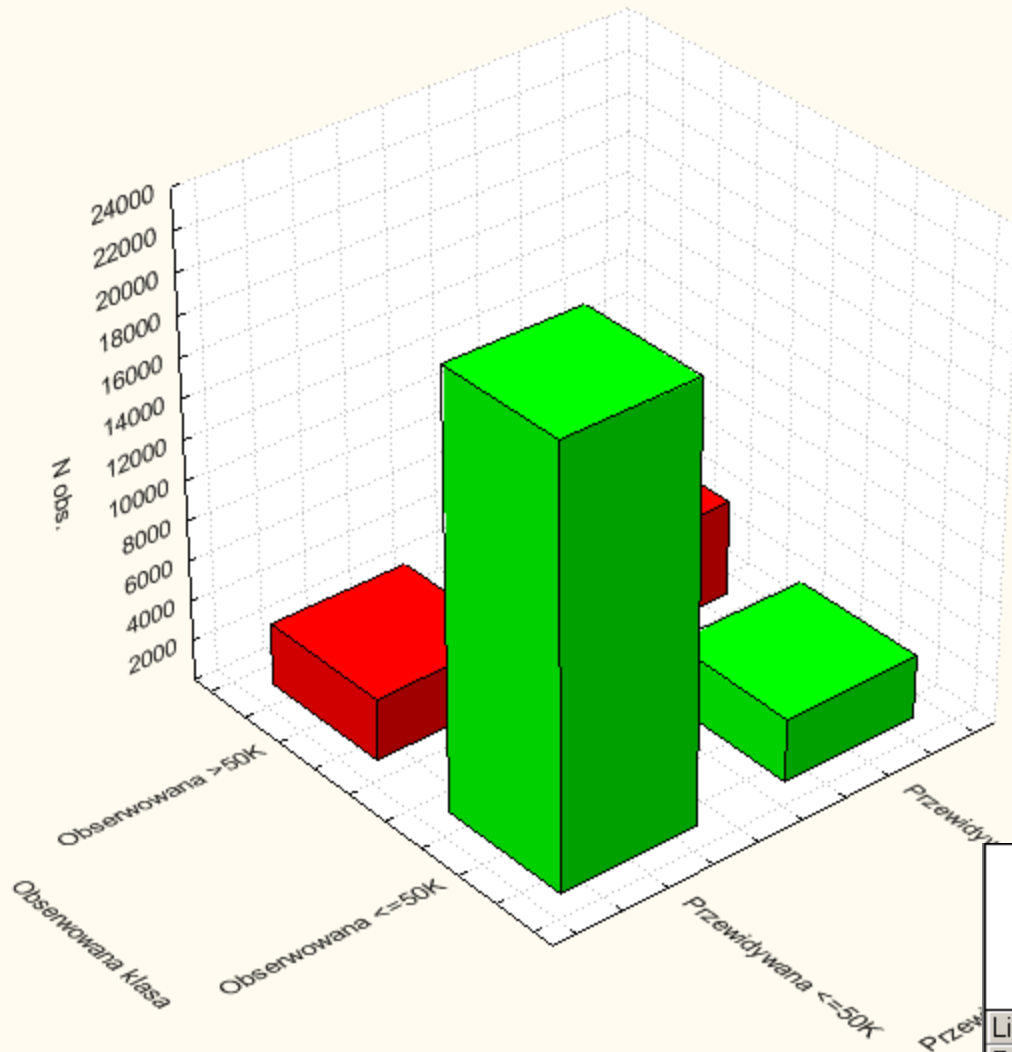


STATISTICA - inny przykład drzewa CHAID





Macierz klasyfikacji



Macierz klasyfikacji (adult)
Zm. zal.: Income
Model: CHAID

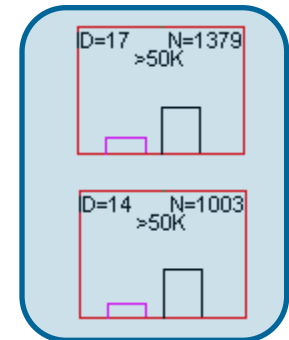
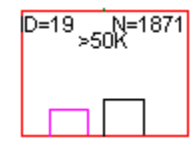
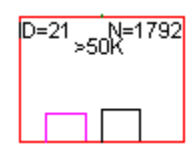
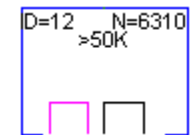
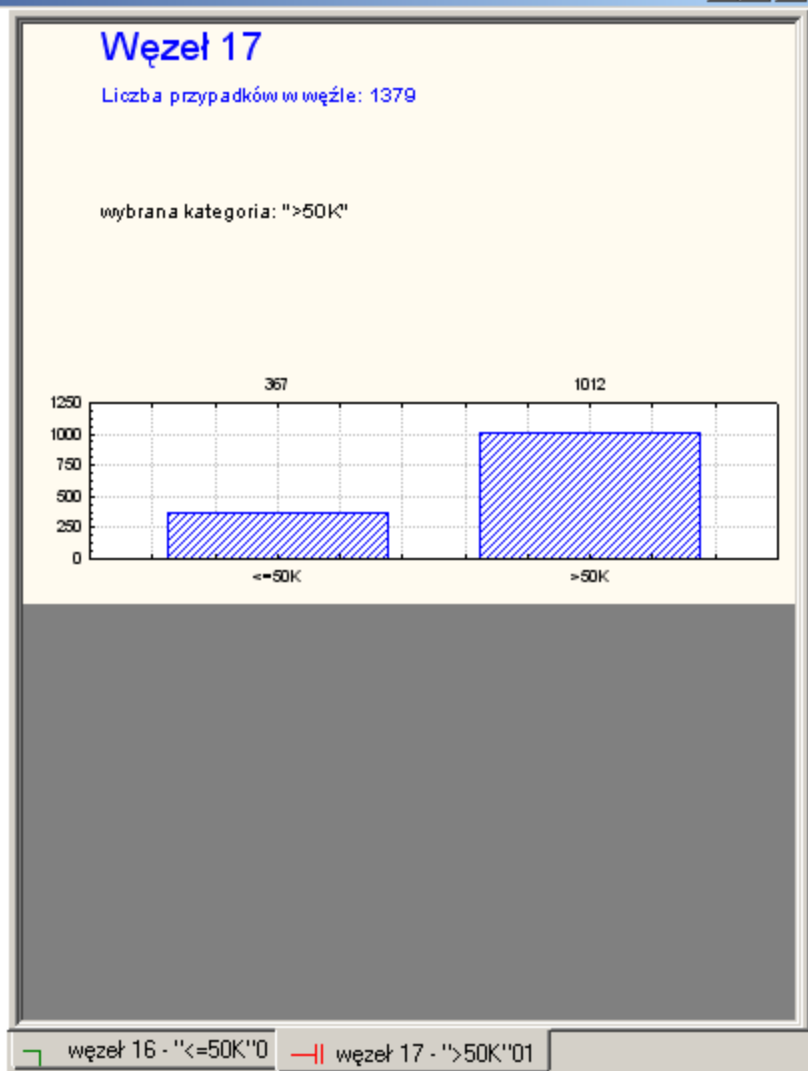
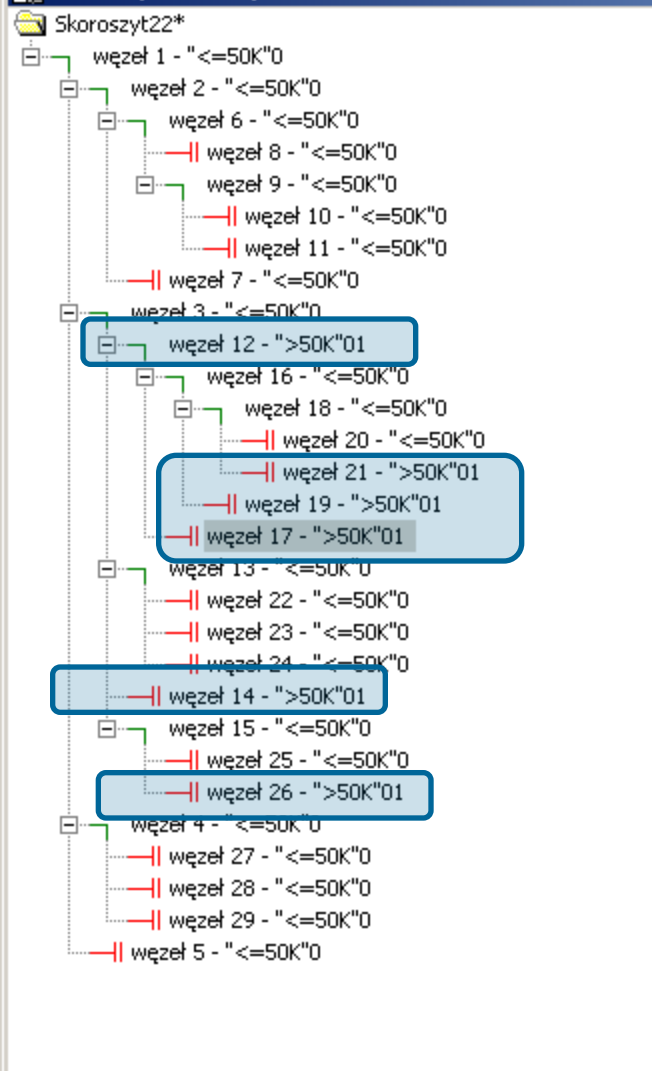
	Obszerw.	Przewidywana <=50K	Przewidywana >50K	Łącznie w wierszu
Liczba	<=50K	21569	3151	24720
Procent z kolumny		87.63%	39.65%	
Procent z wiersza		87.25%	2.75%	
Procent z ogółu		66.24%	9.68%	75.92%
Liczba	>50K	3045	4796	7841
Procent z kolumny		12.37%	60.35%	
Procent z wiersza		88.83%	61.17%	
Procent z ogółu		9.35%	14.73%	24.08%
Liczba	Ogół grup	24614	7947	32561
Procent łącznie		75.59%	24.41%	

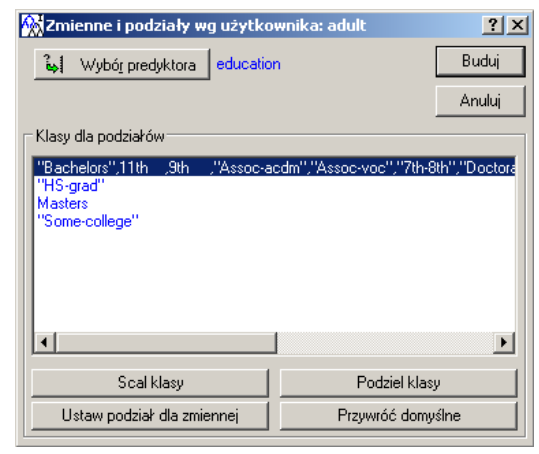
FP

2

FN

1





Jeżeli osoba pozostaje w związku małżeńskim skończyła szkołę z grupy..., ale jest profesjonalistą w swoim zawodzie, wtedy jej dochód prawdopodobnie przekracza 50 000 \$ (węzeł ID17) (z prawdopodobieństwem... 73%)

Jeżeli osoba pozostaje w związku małżeńskim i skończyła studia magisterskie, wtedy jej dochód prawdopodobnie przekracza 50 000 \$ (węzeł ID14) (z prawdopodobieństwem... 77%)

zalety drzew

- łatwe w **interpretacji**, interpretacja w postaci **reguł**
- wyjaśniają **zależności** (*white-box*)
- pozwalają klasyfikować **nieliniowe** problemy
- **szybko** się uczą
- zmienne ilościowe, jak i jakościowe
- nie wymagają założeń co do rozkładu
- pozwalają określić **ważność** predyktorów, niewrażliwe na nieistotne zmienne i wzajemne korelacje
- pełne **pokrycie** przestrzeni wyników
- niewrażliwość na **wartości odstające** – podział w punkcie, nawet jeśli jakieś zmienne osiągają bardzo wysokie/niskie wartości
- radzenie sobie z **brakami danych** – podziały zastępcze

- duża **wariancja** (mała zmiana w zbiorze danych może wywołać dużą zmianę w tworzonych podziałach)
- skłonność do **przeuczania** (bez przycinania drzewa uczą się złożonych, przypadkowych wzorców)
- wymagają wyczucia w ustalaniu **kosztu** dla FN i FP
- *bias* w przypadku nierównej częstości klas
- brak zdolności do **ekstrapolacji** (regresyjnej)
- podejście zachłanne, brak mechanizmów szukania **optimumów globalnych** w szukaniu podziałów

- uczenie zespołowe drzew: Multiple classifiers

Homogeneous classifiers – ten sam algorytm, wiele modeli

- » Bagging (Breiman): Bootstrap aggregation
- » Boosting (Freund, Schapire): AdaBoost, changing the distribution of training examples
- » Multiple partitioned data
- » Multi-class specialized systems, (e.g. ECOC pairwise classification)

Bagging (Bootstrap Aggregation)

- celem jest zmniejszenie wariacji drzew
- budowane jest kilka modeli na różnych podzbiorach (losowanie bootstrapowe)
- niezależne drzewa dokonują predykcji, a wynik jest uśredniany
- przykład: **Random Forest** (losowy las)
- w przypadku losowego lasu losowane są nie tylko przypadki, ale również zmienne

Boosting (drzewa wzmacniane)

- również tworzy zbiór modeli
- podejście sekwencyjne: tworzone są proste drzewa, każde kolejne ma na celu zmniejszać błąd poprzednich
- w każdym kroku rosną wagi przypadków, które zostały błędnie sklasyfikowane
- w ten sposób najtrudniejsze przypadki uzyskują największe znaczenie, a kolejne drzewa zwiększają precyzję
- przykłady [wzmacnianych drzew](#): Gradient Boost, XGBoost, AdaBoost, etc.

single



complete training set

1 iteration



train & keep



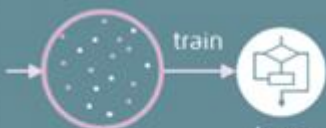
single estimate

reduces over-fitting

bagging



parallel



train & keep



simple average

$$e = \frac{1}{N} \sum_{i=1}^N e_i$$

classification stage

boosting

reduces bias



sequential



train & evaluate



weighted average

$$e = \sum_{i=1}^N w e_i$$