
Eksploracja danych medycznych możliwości wsparcia informatycznego



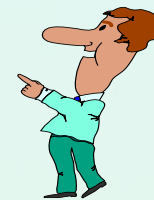
JERZY STEFANOWSKI

Instytut Informatyki
Politechnika Poznańska

ZIM UM 2011

Zawartość

1. Motywacje – dlaczego nowe metody?
2. Rola metod statystycznych
3. Eksploracja danych a systemy uczące
4. Drzewa decyzyjne
5. Indukcja reguł
6. Wybrane zastosowania teorii zbiorów przybliżonych
7. Podsumowanie



1. Motywacje

Zadanie w medycynie → gromadzenie, analiza oraz interpretacja danych



Nowe techniki → „powódź” dostępnych danych

- Urządzenia diagnostyczne
- Rozwój elektronicznych systemów (HSI, medical records ,...)



Rosnące wymagania wobec systemów e-health

Zbyt duża liczba atrybutów / czynników do rozważenia w podejmowaniu decyzji

Dostępne dane mogą mieć różne znaczenie

Dane mogą być niedokładne, nieprecyzyjne i niepełne

Typowe zadanie w analizie medycznych danych:

- Identyfikacja najważniejszych czynników (atrybutów / cech) dla oceny stanu pacjenta
- Odkrywanie zależności między opisem pacjenta (wartości atrybutów) a decyzją co do pacjenta, np. klasyfikacją pacjentów (diagnozy, sposoby postępowania)
- Inne

Typowe obszary zastosowań

Za R. E. Abdel-Aal „Data mining and medical informatics”

- Screening
- Diagnosis
- Therapy
- Prognosis
- Monitoring
- Biomedical/Biological Analysis
- Epidemiological Studies
- Hospital Management
- Medical Instruction and Training

Analiza ważności atrybutów

- Dane – wielowymiarowe tablice (pacjenci / atrybuty)
- Wiele możliwych podejść:
 - Statystyczna analiza danych
 - Statystyczne miary siły związku.
 - Metody wielowymiarowe (analiza czynnikowa, skalowanie wielowymiarowe, analiza dyskryminacyjna, ...)
 - Metody wywodzące się z Uczenia Maszynowego (Machine Learning) i Data Mining (Eksploracja danych)

Statystyka opisowa a wnioskowanie statystyczne

- **Statystyka opisowa** – prezentacja danych w sposób uporządkowany, prosty z wykorzystaniem, np., miar tendencji centralnej i miar rozproszenia. Pomagają one w redukcji dużej liczby danych do zbioru bardziej zwężonych i "pojemnych" miar.
- **Wnioskowanie statystyczne** – pozwala ustalać prawidłowości i podejmować decyzje na podstawie zredukowanej liczby danych (próby) przy zastosowaniu rachunku prawdopodobieństwa.
- Rachunek prawdopodobieństwa – możliwe jest określenie jak błąd popełnia się uogólniając wyniki z **próby** na całą **zbiorowość**

2. Weryfikacja hipotez statystycznych

- Testy parametryczne → **sprawdzenie pewnej hipotezy** dotyczącej poziomu nie znanego parametru albo co do postaci rozkładu zmiennej w populacji.
- Na podstawie informacji pochodzącej **z próby** będziemy podejmować decyzje czy przyjąć albo odrzucić hipotezę.
- Przykłady problemów badawczych dotyczących:
 - wartości badanych zmiennych,
 - np. średni wiek osób chorujących na pewną chorobę wynosi 45 lat.
 - porównania dwóch zbiorowości,
 - skuteczność oddziaływania pewnych bodźców, którym poddawane są te same grupy obiektów,
 - zależności między badanymi zmiennymi,
 - porównania rozkładów zmiennych.

Przykład testu dla 2 zbiorowości (t-Studenta)

Microsoft Excel - T-TEST2.XLS

Wpisz pytanie do Pomocy

Plik Edycja Widok Wstaw Format Narzędzia Dane Okno Pomoc Adobe PDF

100% Arial

K14 =TEST.T(B10:B31;D10:D29;2;2)

zakłada się równe wariancej - Patrz Oktaba rozdział 7.5
 Uwaga są drobne różnice w obliczeniach w stosunku do przykładu z książki Pana Oktaby
 Ale wniosek co do odrzucenia hipotezy zerowej H_0 jest ten sam - można ją odrzucić!

Przykład 7.5.2
 Wysokość czaszek ryjówek w miesiącach czerwcu i lipcu

	lipiec	czerwiec		
	6,60	6,5		
	6,60	6,5		
	6,50	6,4		
	6,50	6,4		
	6,40	6,3		
	6,40	6,2		
	6,40	6,2		
	6,40	6,1		
	6,40	6,1		
	6,40	6,1		
	6,40	6,1		
	6,40	6,1		
	6,40	6,1		
	6,40	6,1		
	6,30	6,1		
	6,30	6,1		
	6,30	6,00		
	6,30	6,00		
	6,30	6,00		
	6,30	6,00		
	6,30	6,00		
	6,30	6,00		
	6,20	5,90		
	6,10	5,90		
	6,10			
	6,00			
licznosc	22,00	20,00		
srednia	6,34	6,15		

Efekt uzycia opcji z Data Analysis wykorzystanie dostępnych funkcji

t-Test: Two-Sample Assuming Equal Variances

	czerwiec	lipiec	
Mean	6,34	6,15	t(0.05) 2,02
Variance	0,02	0,03	p(T0) 0,000505
Observations	22,00	20,00	
Pooled Variance	0,03		
Hypothesized Mean Difference	0,00		
df	40,00		stopnie swobody
t Stat	3,78		wartość wyliczona statystyki t -Studenta
P(T<=t) one-tail	0,00		test jednostronny
t Critical one-tail	1,68		wartość krytyczna dla 5%
P(T<=t) two-tail	0,00		test dwustronny
t Critical two-tail	2,02		wartość krytyczna dla 5%

Wniosek Hipotezę można odrzucić

t_Student 1 | t-Stud ver 1 | t-Stud ver 2 | t_Student TwoPaired | t_Student TwoPaired

Gotowy

Test niezależności zmiennych - Chi-kwadrat

STATISTICA: Podstawowe statystyki i tabele

Plik Edycja Widok Analiza Wykresy Opcje Okno Pomoc

1. Kolumny Wiersze

Dane: chikwadrat.sta 4v * 40c

Test chi kwadrat obliczenia kontrolne

	1 SKALA	2 PAPIEROS	3 ALKOHOL	4 PRACA
1	III	duzo	duzo	duzo
2	I	duzo	malo	duzo
3	III	duzo	duzo	duzo
4	III	duzo	duzo	duzo
5	III	duzo	malo	duzo
6	I	malo	duzo	srednio
7	III	srednio	duzo	duzo
8	II	duzo	duzo	malo
9	III	malo	malo	nic
10	II	duzo	duzo	nic
11	III	duzo	srednio	duzo

Tabela liczebności (chikwadrat.sta)

PODST. Licznosc oznacz. komorek > 5
STATYST. (Nieoznaczono sum brzegowych)

PAPIEROS	ALKOHOL nic	ALKOHOL malo	ALKOHOL srednio	ALKOHOL duzo	Wiersz Razem
malo	1	4	1	2	8
srednio	2	4	2	1	9
duzo	2	6	7	8	23
Ogól grp	5	14	10	11	40

Zestawienie: Licznosci oczekiwane (chikwadrat.sta)

PODST. Licznosc oznacz. komorek > 5
STATYST. Chi^2 Pearsona: 4,35288, df=6, p=.629039

PAPIEROS	ALKOHOL nic	ALKOHOL malo	ALKOHOL srednio	ALKOHOL duzo	Wiersz Razem
malo	1,000000	2,800000	2,000000	2,200000	8,000000
srednio	1,125000	3,150000	2,250000	2,475000	9,000000
duzo	2,875000	8,050000	5,750000	6,325000	23,000000
Ogól grp	5,000000	14,000000	10,000000	11,000000	40,000000

Wynik tabelaryzacji

Przeгляд tabeli zbiorczej

Dokładne tabele dwudzielcze

Tabela zbiorcza

Pokaż długie etykiety wartości

Wliczaj braki danych

Pokaż wybrane % w oddzielnych tabelach

Statystyki dla tabel dwudzielczych

Chi-kwadrat Pearsona i NW

dokładny Fishera, Yatesa, McNemara (2 x 2)

Fi (tabela 2x2), V i C Craméra

tau-b i tau-c Kendalla

Gamma

Współczynnik korelacji rang Spearmana

d Sommera

Współczynniki niepewności

Tabele

Podświetl liczebności > 5

Liczebności oczekiwane

Liczebności resztowe

Procenty z całości

Procenty w wierszach

Procenty w kolumnach

Skategoryzowane histogramy

Wykresy interakcji liczebności

Histogramy 3W

UWAGA: Tabele zbiorcze są tworzone wyłącznie jeżeli zostały wybrane dwie listy zmiennych. Aby wyznaczyć Chi-kwadrat największej wiarygodności i analizować wielodzielcze tabele liczebności używamy analizy logliniowej.

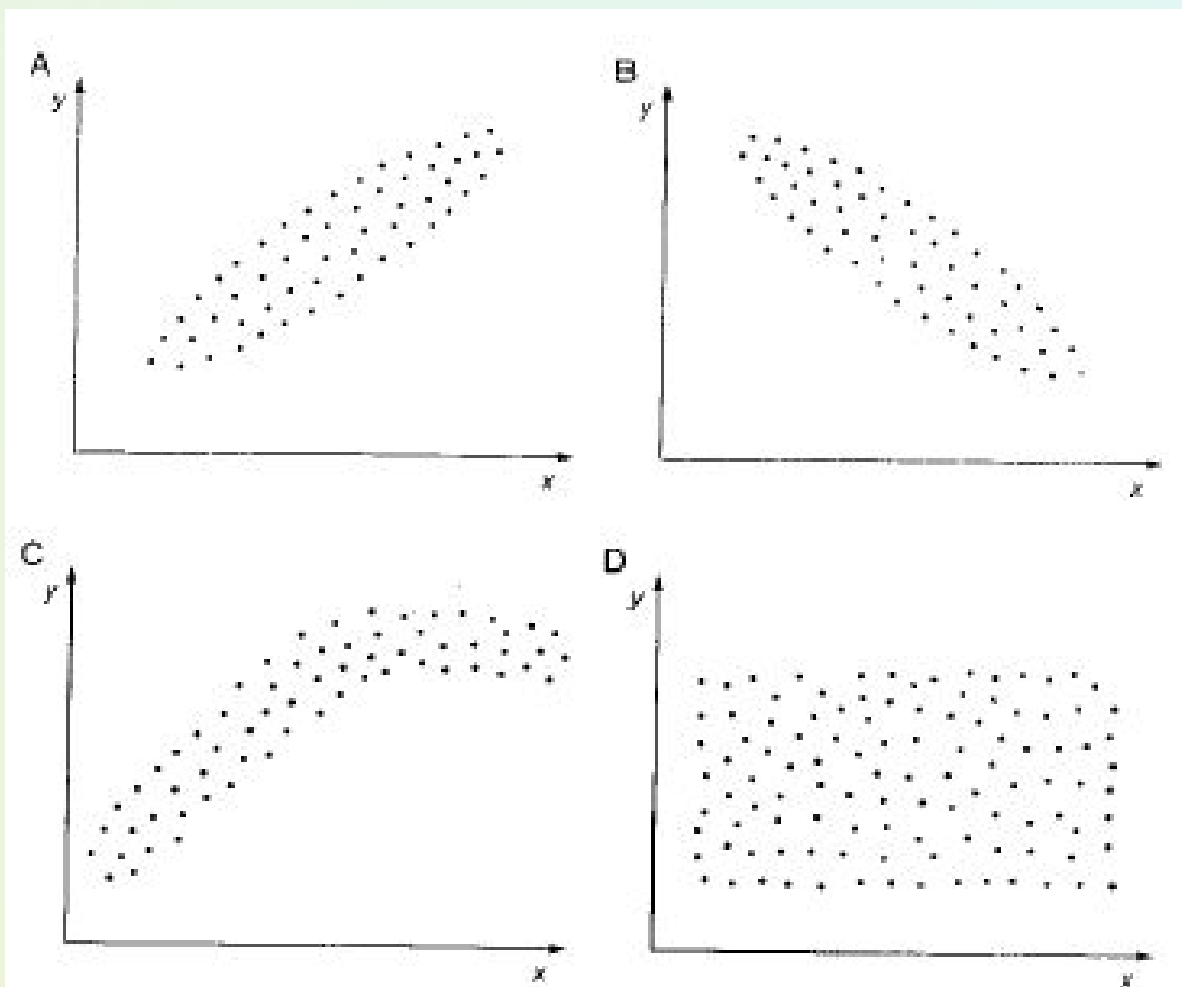
1: Rozkład dwuwymiarowy: PAPIEROS x ALKOHOL

Rozkład dwuwymiarowy: PAPIEROS x ALKOHOL

Wyjście: WYŁĄCZONE | Sel: NIE | Waga: WYŁĄCZONA

Start | Menedżer zada... | 2 Microsoft P... | Total Comman... | STATISTICA: P... | Przełącznik mo... | PL | 11:33

Różne zależności między zmiennymi



Rys. 10.1. Rozkłady empiryczne dwóch zmiennych



Analiza korelacji (liniowej) – pakiet Statistica

STATISTICA - Workbook2* - [Correlations (EnginePerformance.sta)]

File Edit View Insert Format Statistics Graphs Tools Data Workbook Window Help

Clipboard Paste Paste as Text Paste as Text with Formatting Paste as Text with Merged Cells Paste as Text with Row and Column Spacing Add to Workbook Add to Report

Arial 10 B I U

Data: EnginePerformance.sta (79v by 128c)

	1	2	3	4	5	6	7	8	9	10	11
	Serial Number	Efficiency	Fuel Economy(%)	Power(%)	Input01	Input02	Input03	Input04	Input05	Input06	Input07
1	#25457	102,384	100,066	99,814	100,186545	16,6255147	11,9297997	15,4501075	16,7199319	17,4754064	20,753
2	#25458	81,405	89,798	110,392	98,4136317	16,3445083	13,5326772	14,0013087	15,6347214	17,050197	20,303
3	#25459	94,070	92,072	87,917	98,7403916	16,5964348	12,0007502	15,5077475	15,7857113	18,6175749	20,527
4	#25460	108,855	89,369	90,945	99,5529412	16,7615965	12,0610633	14,2580726	13,8695801	17,8851961	19,81
5	#25461	107,903	89,453	95,912	98,8236109	16,6525248	12,2789147	14,6501313	20,634384	17,1218605	21,11
6	#25462	86,475	94,063								
7	#25463	105,583	94,868								
8	#25464	109,303	95,652								
9	#25465	103,633	91,181								
10	#25466	95,300	93,490								
11	#25467	102,334	90,320								
12	#25468	94,456	118,944								
13	#25469	109,349	107,966	1							
14	#25470	105,943	89,392								
15	#25471	101,390	102,309								
16	#25472	105,911	107,008	1							
17	#25473	78,027	91,527								
18	#25474	107,266	89,611								
19	#25475	99,571	101,998	1							
20	#25476	107,466	102,613	1							
21	#25477	109,327	95,364	1							
22	#25478	104,091	91,369								
23	#25479	95,655	90,542								
24	#25480	107,033	96,745								
25	#25481	108,802	107,768	1							
26	#25482	98,975	117,309	1							
27	#25483	104,152	100,064	1							
28	#25484	67,792	116,900								

Workbook2* - Correlations (EnginePerformance.sta)

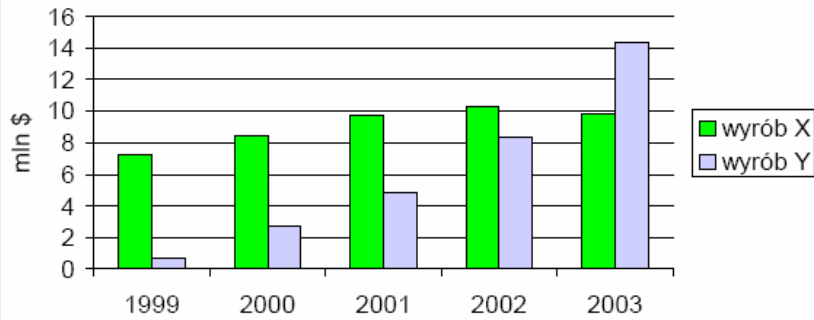
Correlations (EnginePerformance.sta)
Marked correlations are significant at $p < ,05000$
N=128 (Casewise deletion of missing data)

Variable	Efficiency	Fuel Economy(%)	Power(%)	Input01	Input02	Input03
Efficiency	1,00	-0,09	0,12	0,12	0,19	0,
Fuel Economy(%)	-0,09	1,00	0,53	0,67	0,50	0,
Power(%)	0,12	0,53	1,00	0,26	0,14	0,
Input01	0,12	0,67	0,26	1,00	0,83	-0,
Input02	0,19	0,50	0,14	0,83	1,00	-0,
Input03	0,06	0,10	0,12	-0,01	-0,05	1,
Input04	-0,07	-0,08	0,00	-0,20	-0,23	-0,
Input05	-0,00	-0,00	0,06	-0,10	-0,04	0,
Input06	0,15	0,11	0,17	0,14	0,16	0,

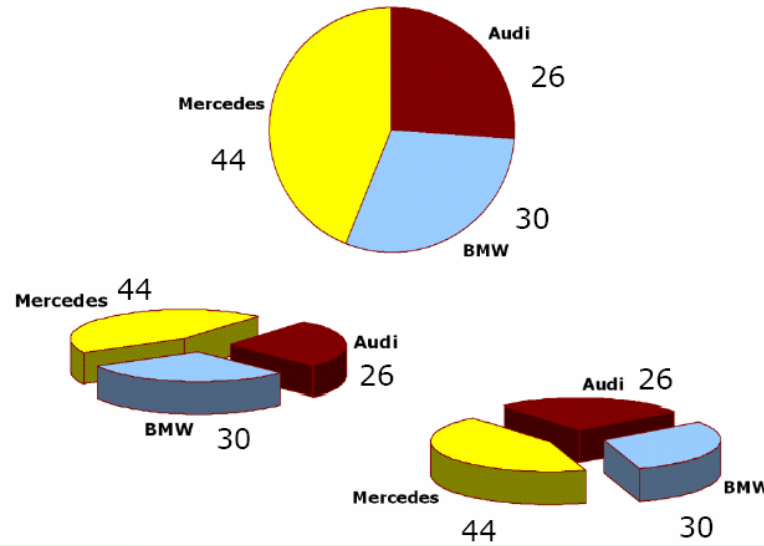
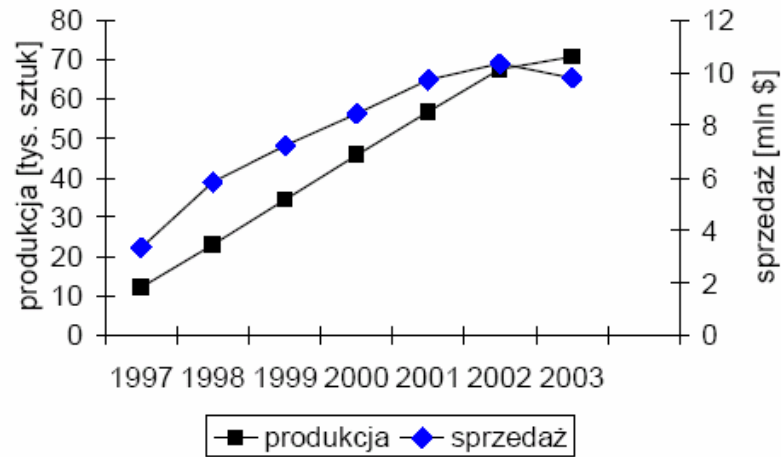
Wizualizacja informacji

- Elementy graficzne → tabele, wykresy, diagramy, rysunki, fotografie, mapy
- UZASADNIENIE
 - ludzie są „wzrokowcami”: 83% przyswajanej wiedzy wynika z pobudzeń wzrokowych
 - człowiek zapamiętuje ok. 43% więcej informacji, jeśli dokument zawiera elementy graficzne
- RÓŻNE FUNKCJE
 - przekazywanie informacji, która jest trudna do przekazania w inny sposób (słowami)
 - pomoc w wyjaśnianiu i wyróżnianiu informacji

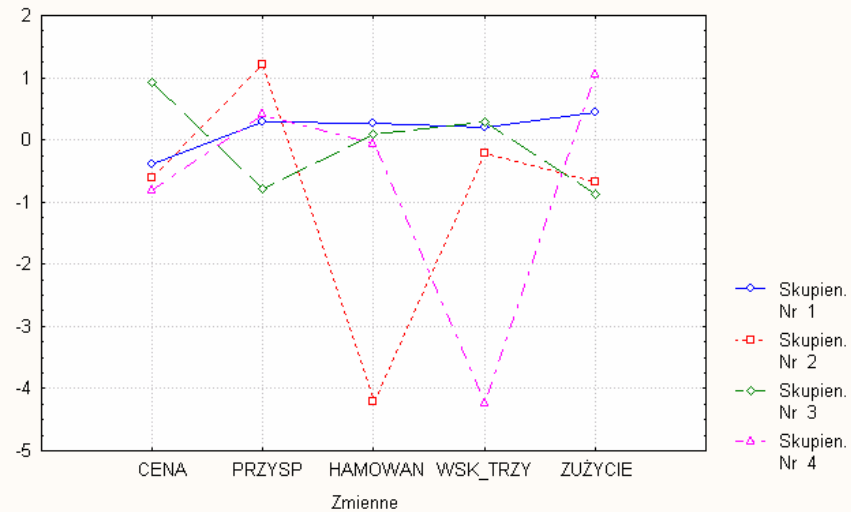
sprzedaż w latach 1999-2003



wyrób X - produkcja i sprzedaż w latach 1997-2003



Wykres średnich każdego skupienia



- Przykłady różnych wykresów

Inne ciekawe metody

- Statystyka opisowa
 - Miary dynamik zjawisk; metody indeksowe
- Estymacja przedziałowa; przedziały ufności
- Analiza wariancji – ANOVA
- Regresja wielokrotna (także inne zagadnienia)
- Metody nieparametryczne i inne metody analizy danych jakościowych
- Kontrola i poprawa jakości
- Analiza wielowymiarowa
 - Wielowymiarowa analiza wariancji, analiza dyskryminacyjna, metoda głównych składowych i analiza czynnikowa, skalowanie wielowymiarowa
 - Analiza dyskryminacyjna i metody klasyfikowania
 - Analiza skupień
- Model bayesowski i analiza decyzji
- Metody wyboru prób z populacji
- ...



No i jeszcze trochę o software

MINITAB [Statistical Software](#) jest archetypem wśród programów statystycznych. Jego pierwsza wersja powstała przed 25 laty. Od 1972 roku był wykorzystywany i cytowany w ponad 300 publikacjach. Jest programem prostym w obsłudze i łatwym do nauczenia.

STATISTICA [StatSoft](#) - oferuje bogaty zestaw metod statystycznych i operacji na danych, jest łatwy w obsłudze, ma znakomitą grafikę i narzędzia do tworzenia interaktywnych aplikacji. Lista dostępnych funkcji jest bardzo długa, w jej skład wchodzi między innymi: podstawowe i zaawansowane wielowymiarowe metody statystyczne. Także uproszczone karty kontrolne, analiza przeżycia, sterowanie jakością, analiza procesu, planowanie doświadczeń oraz sieci neuronowe. Wersja Data Miner – metody data mining!

SPSS, SPSS Inc. - bardzo silne zaawansowane narzędzie. Podstawowym elementem pakietu jest moduł SPSS Base zawierający podstawowe statystyki, wielowymiarowe tablice i wykresy. Dodatkowe moduły to Advanced Statistics - zaawansowane funkcje.

SAS, SAS Institute – historycznie podstawowy zestaw modułów przetwarzania danych dla komputerów „mainframe”. Wiele procedur statystycznych. Wsparcie do integracji danych i konstruowania hurtowni danych oraz OLAP.



No i jeszcze dalej o software

STATGRAPHICS Plus [Statistical Graphics Corporation](#) - jeden z najpopularniejszych pakietów, jest idealnym narzędziem do celów dydaktycznych oraz do analizy danych o małych i średnich rozmiarach

MATLAB – Statistics Toolbox. Oryginalnie oprogramowanie The MathWorks, Inc. - przeznaczony głównie do realizacji obliczeń numerycznych. Uzupełniony w 1991 roku pakietem SIMULINK, który między innymi przejął funkcję interfejsu, wprowadzając elementy graficznych obiektów oraz okienkowy tryb dialogu z użytkownikiem

S-PLUS [MathSoft](#) - przeznaczony do prowadzenia wszechstronnej analizy statystycznej danych i generowania prezentacji graficznych. Zawiera wszystkie opcje języka S, który jest specjalizowany do programowania kompleksowej analizy statystycznej i wizualizacji danych

SYSTAT - kolejny pakiet [SPSS, Inc.](#) Posiada wiele funkcji przetwarzania danych i prezentacji wyników. Zawiera klasyczne procedury statystyk opisowych, testów nieparametrycznych, korelacji, szeregów czasowych, dostępne są modele liniowe i logliniowe, analizy - skupień, przeżycia, analiza dyskryminacyjna i inne

Sigma Plot 2000 [Jandel Scientific Software](#) - program przetwarzania danych, posiada rozbudowaną opcję kreślenia wykresów, ogromne możliwości edytorskie spowodowały, że pakiet wykorzystywany jest do przygotowywania profesjonalnych wizualizacji graficznych



R project

The screenshot displays the R GUI interface with several windows open:

- R Console:** Shows the following R code:

```
rgl.sr> ylen <- ylim[2] - ylim[1] + 1
rgl.sr> colorlut <- terrain.colors(ylen)
rgl.sr> col <- colorlut[y - ylim[1] + 1]
rgl.sr> rgl.clear()
rgl.sr> rgl.surface(x, z, y, color = col)
```
- R Data Editor:** Displays a table with columns 'height' and 'weight':

height	weight
58	115
59	117
60	120
61	123
62	126
63	129
64	132
65	135
66	139
67	142
68	146
69	150
70	154
71	159
72	164
- Quartz (2) - Active:** Shows a 2D plot titled 'Given : depth' with 'long' on the x-axis and 'depth' on the y-axis. The plot displays a series of points forming a curve.
- R Workspace Browser:** Lists objects in the workspace:

Object	Type	Structure
dati	data.frame	dim: 20 4
g	factor	levels: 10
l	numeric	length: 12
n	numeric	length: 1
opar	list	length: 2
pie.sales	numeric	length: 6
pin	numeric	length: 2
scale	numeric	length: 1
usr	numeric	length: 4
women	data.frame	dim: 15 2
height	numeric	length: 15
weight	numeric	length: 15
x	numeric	length: 87
- R Package Manager:** Lists installed and available packages:

status	Package	Description
<input checked="" type="checkbox"/> loaded	graphics	The R Graphics Package
<input type="checkbox"/> not loaded	grid	The Grid Graphics Package
<input type="checkbox"/> not loaded	lattice	Lattice Graphics
<input checked="" type="checkbox"/> loaded	methods	Formal Methods and Classes
<input type="checkbox"/> not loaded	mgcv	GAMs with CCV, smoothness estimation
- RGL device 1 (active):** Shows a 3D surface plot of a terrain, colored with a gradient from green to yellow to red.

R project



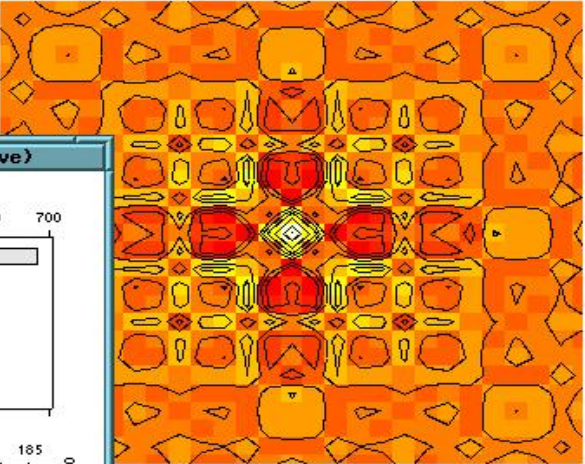
```
leisch@galadriel:~/work/tnp
R> n <- 5
R> g <- gl(n, 100, n*100)
R> x <- rnorm(n*100) + sqrt(codes(g))
R> boxplot(split(x,g), col="lavender", notch=TRUE)
R> title(main="Notched Boxplots", xlab="Group", font.main=4, font.lab=1)
R>
R> ctl <- c(4.17,5.58,5.18,6.11,4.50,4.61,5.17,4.53,5.33,5.14)
R> trt <- c(4.81,4.17,4.41,3.59,5.87,3.83,6.03,4.89,4.32,4.69)
R> group <- gl(2,10,20,labels=c("Ctl","Trt"))
R> weight <- c(ctl,trt)
R> anova(lm,D9 <- lm(weight~group))

Analysis of Variance Table
Response: weight
          Df Sum Sq Mean Sq    F Pr(>F)
group     1  0.6882   0.6882  1.419  0.249
Residual 18  8.7293   0.4850

R>
R> 
```

R Graphics: Device 2 (inactive)

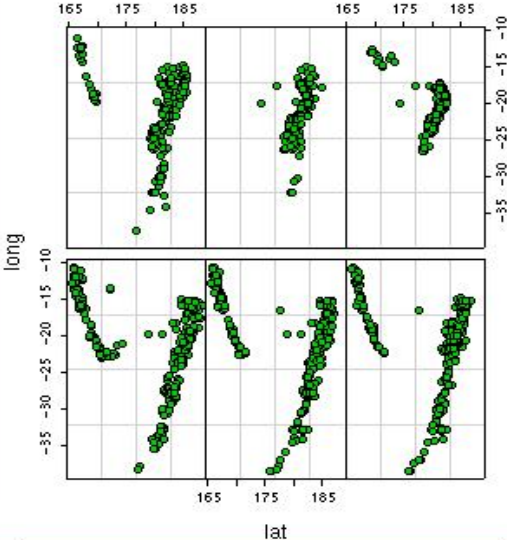
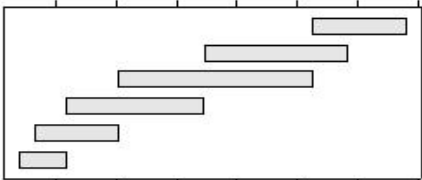
Math can be beautiful ...



$\cos(r^2)e^{-r^{1.6}}$

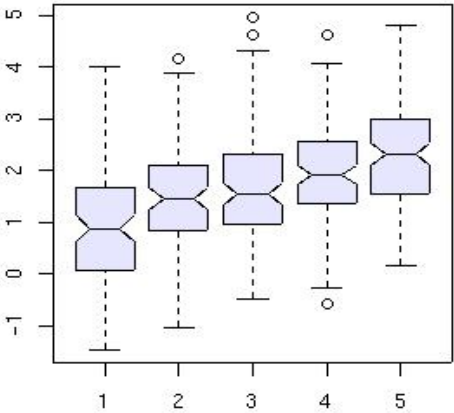
R Graphics: Device 3 (inactive)

Given : depth



R Graphics: Device 4 (ACTIVE)

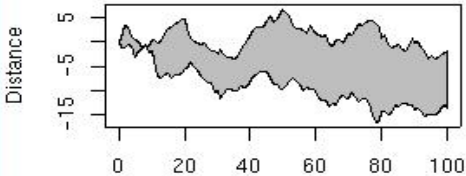
Notched Boxplots



Group

R Graphics: Device 5 (inactive)

Distance Between Brownian Motions



Distance

Time

Dlaczego poszukiwać nowych metod?

Za R. E. Abdel-Aal „Data mining and medical informatics”

- Manual hypothesis testing:

Not practical with large numbers of variables

- User-driven... User specifies variables, functional form and type of interaction:

User intervention may influence resulting models

- Assumptions on linearity, probability distribution, etc.

May not be valid

- Datasets collected with statistical analysis in mind

Not always the case in practice

Rozwój środków obliczeniowych

- Tańsze, pojemniejsze i szybsze pamięci
- Technologicznie można gromadzić TB danych
- Dojrzały stan systemów zarządzania bazami danych
- Rozwój metod (automatycznego) pozyskiwania danych
- Coraz większą moc obliczeniowa

Sztuczna inteligencja

- SI. jest nauk o maszynach realizujących zadania, które wymagają inteligencji wówczas, gdy są wykonywane przez człowieka (Minski)
- SI stanowi dziedzinę informatyki dotyczącą metod i technik wnioskowania symbolicznego przez komputer oraz symbolicznej reprezentacji wiedzy stosowanej podczas takiego wnioskowania (Feigenbaum).
- Sztuczna Inteligencja, to dziedzina informatyki zajmująca się rozwiązywaniem zadań efektywnie „niealgorytmizowalnych” w oparciu o modelowanie wiedzy.



Typowe działy sztucznej inteligencji

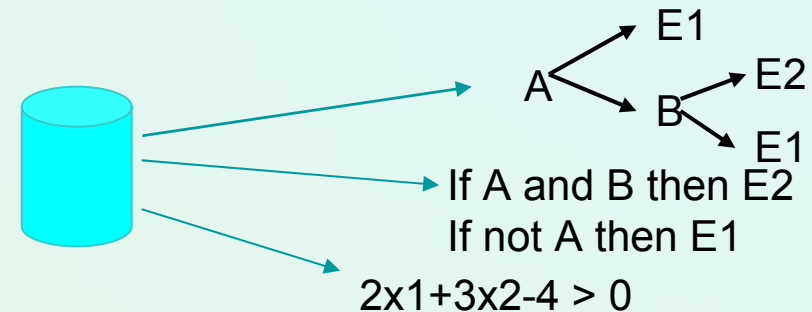
- Programy i maszyny grające (teoria gier)
- Systemy eksperckie (ekspertowe)
- Rozpoznawanie języka naturalnego
- Procesy percepcji
- Uczenie maszynowe
- Inteligentne przeszukiwanie danych
- Robotyka
- tzw. sztuczne życie
- Dowodzenie twierdzeń matematycznych

Zaawansowane dziedziny informatyki

- **Uczenie maszynowe**
→ (ang. **Machine Learning**)
- Systemy, które doskonałą swoje działanie na podstawie

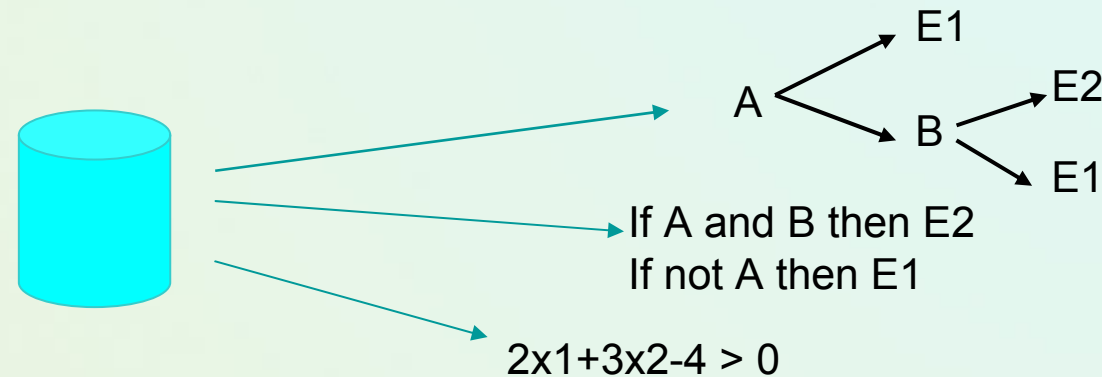


DARPA LAGR Herminator



- **Eksploracja danych** → **Data mining**
 - Poszukiwanie w zgromadzonych danych nieznanymi, użytecznymi regularnościami, związkami między elementami danych.
 - Potencjalnie duże / złożone repozytoria danych
- Odkrywanie wiedzy.

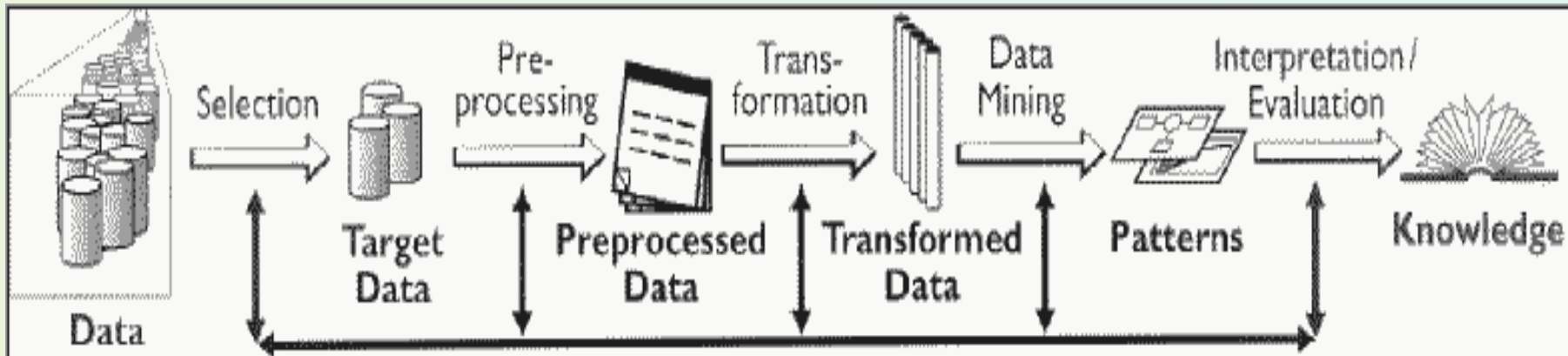
Eksploracja danych



- **Data mining** (ang.)
 - Poszukiwanie w zgromadzonych danych nieznanymi, użytecznymi regularnościami, związków między elementami danych.\
- Eksploracja danych to **etap** w procesie odkrywania wiedzy.
- **Wiedza** → uporządkowana i formalna reprezentacja odkrytej regularności między elementami danych

Odkrywanie Wiedzy w Bazach Danych

Def: „Nietrywialny proces poszukiwania nowych, użytecznych i zrozumiałych wzorców (regularności) z danych” [Piatetsky]



Przetwarzanie wstępne:

Selekcja – wybór podzbioru przykładów i atrybutów

Pre-processing – wstępne przetwarzanie, czyszczenie danych, standaryzacja

Transformacje – do postaci akceptowanej przez metody eksploracji danych

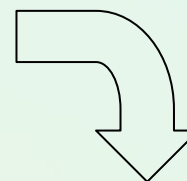
Wiedza klasyfikacyjna

- Problem określania zasad przydziału obiektów do znanych wstępnie klas na podstawie analizy danych o przykładach klasyfikacji.

Wiek	Zawód	dochód	...	Decyzja
21	Prac. fiz.	1220	...	Nie kupi
26	Menedżer	2900	...	Kupuje
44	Inżynier	2600	...	Kupuje
23	Student	1100	...	Kupuje
56	Nauczyciel	1700	...	Nie kupi
...
45	Lekarz	2200	...	Nie kupi
25	Student	800	...	Kupuje

Przykłady uczące

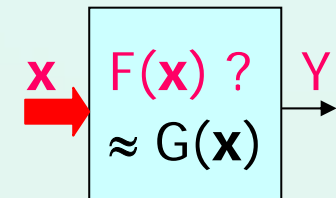
Algorytm eksploracji



Reprezentacja wiedzy:
np. reguły
R1. Jeżeli student to kupuje komputer
R2. Jeżeli dochód > 2400 ...

Klasyfikacja nadzorowana - Supervised Learning

- $Y=F(\mathbf{x})$: true function (usually not known) for population P
- Przykłady etykietowane $\langle \mathbf{x}, Y \rangle$
- 1. Collect Data: “labeled” training sample drawn from P



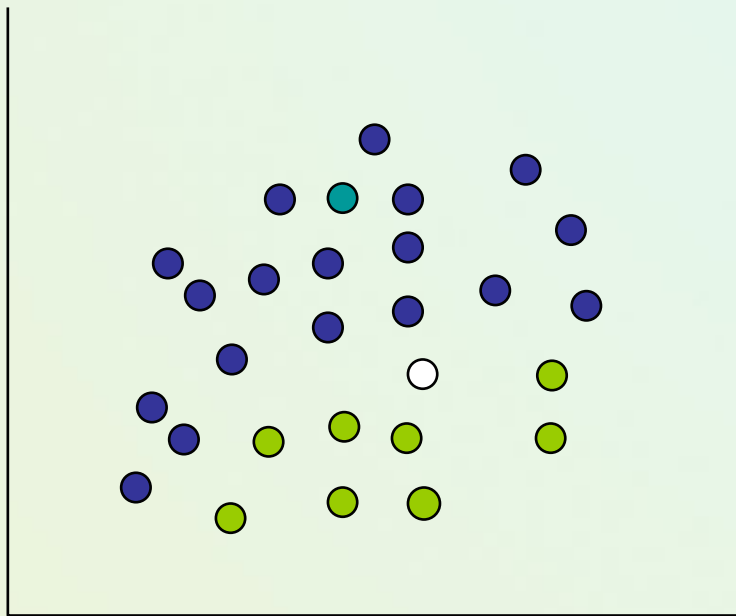
57,M,195,0,125,95,39,25,0,1,0,0,0,1,0,0,0,0,0,0,1,1,0,0,0,0,0,0,0,0 0
78,M,160,1,130,100,37,40,1,0,0,0,1,0,1,1,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0 1
69,F,180,0,115,85,40,22,0,0,0,0,0,1,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,0 0
18,M,165,0,110,80,41,30,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0 1

- 2. Training: Get $G(\mathbf{x})$; model learned from training sample,
Goal: $E\langle (F(\mathbf{x})-G(\mathbf{x}))^2 \rangle \approx 0$ for future samples drawn from P – Not just data fitting!
- 3. Test/Use:

71,M,160,1,130,105,38,20,1,0,0,0,0,0,0,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0 ?

Klasyfikowanie

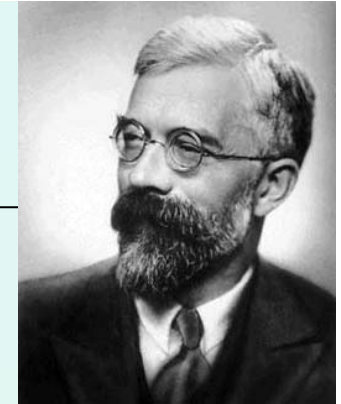
Odkryj z danych historycznych metody przydziału przykładów do klas. Przykłady uczące - są opisane zbiorem atrybutów i etykietowane



Wiele podejść :
Analiza dyskryminacyjna
Klasyfikatory Bayesowskie,
Sieci neuronowe,
Drzewa i reguły decyzyjne,
K-NN
...

Dane przykłady uczące z klas ● ●
Określ przydział dla nowego obiektu ○?

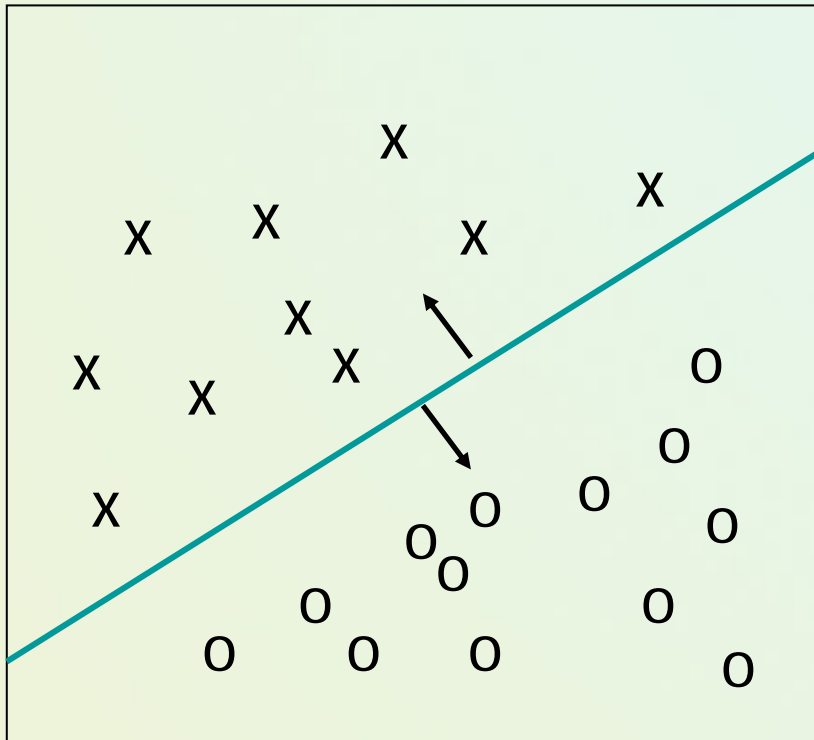
Klasyfikacja – podejście statystyczne



$$D = \left\{ (\mathbf{x}_i, c_i) \mid \mathbf{x}_i \in R^p, c_i \in \{C_1, \dots, C_k\} \right\}_{i=1}^N$$

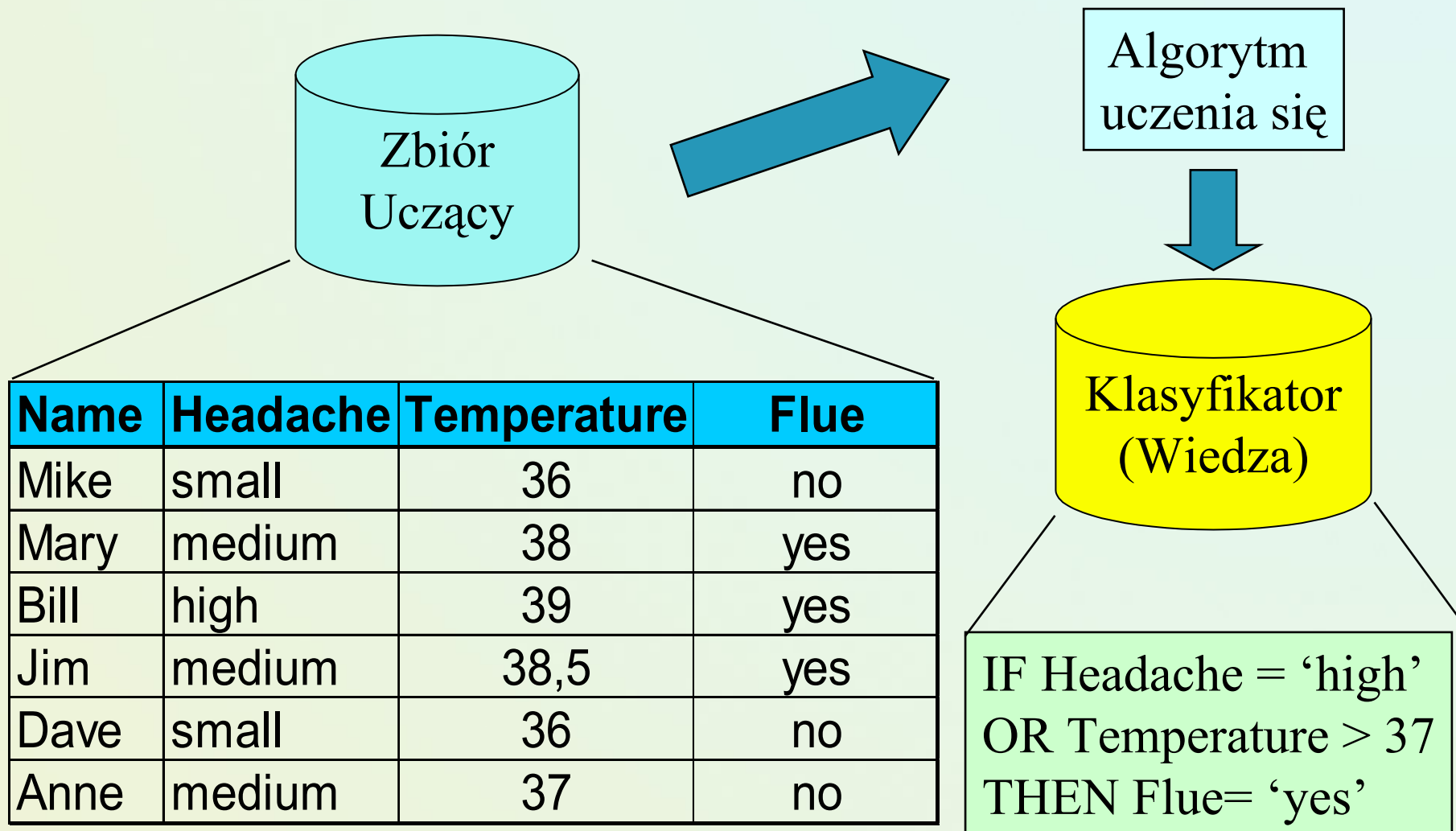
$$y = f(\mathbf{x}, \mathbf{w})$$

$$y = \begin{cases} f(\mathbf{x}_i) > T & \mathbf{x}_i \in C_1 \\ f(\mathbf{x}_i) < T & \mathbf{x}_i \in C_2 \end{cases} \quad \text{sir Ronald Fisher}$$



- Binarna klasyfikacja (uogólnienie na więcej klas)
- Poszukiwanie przybliżenia „granicy decyzyjnej” – ang. decision boundary
- Obserwacje ponad linią przydziel do klasy ‘x’
- Obserwacje pod linią przydziel do klasy ‘o’
- Przykłady: Fisher-owska analiza dyskryminacyjna, SVM, ANN

Poszukiwanie wiedzy klasyfikacyjnej



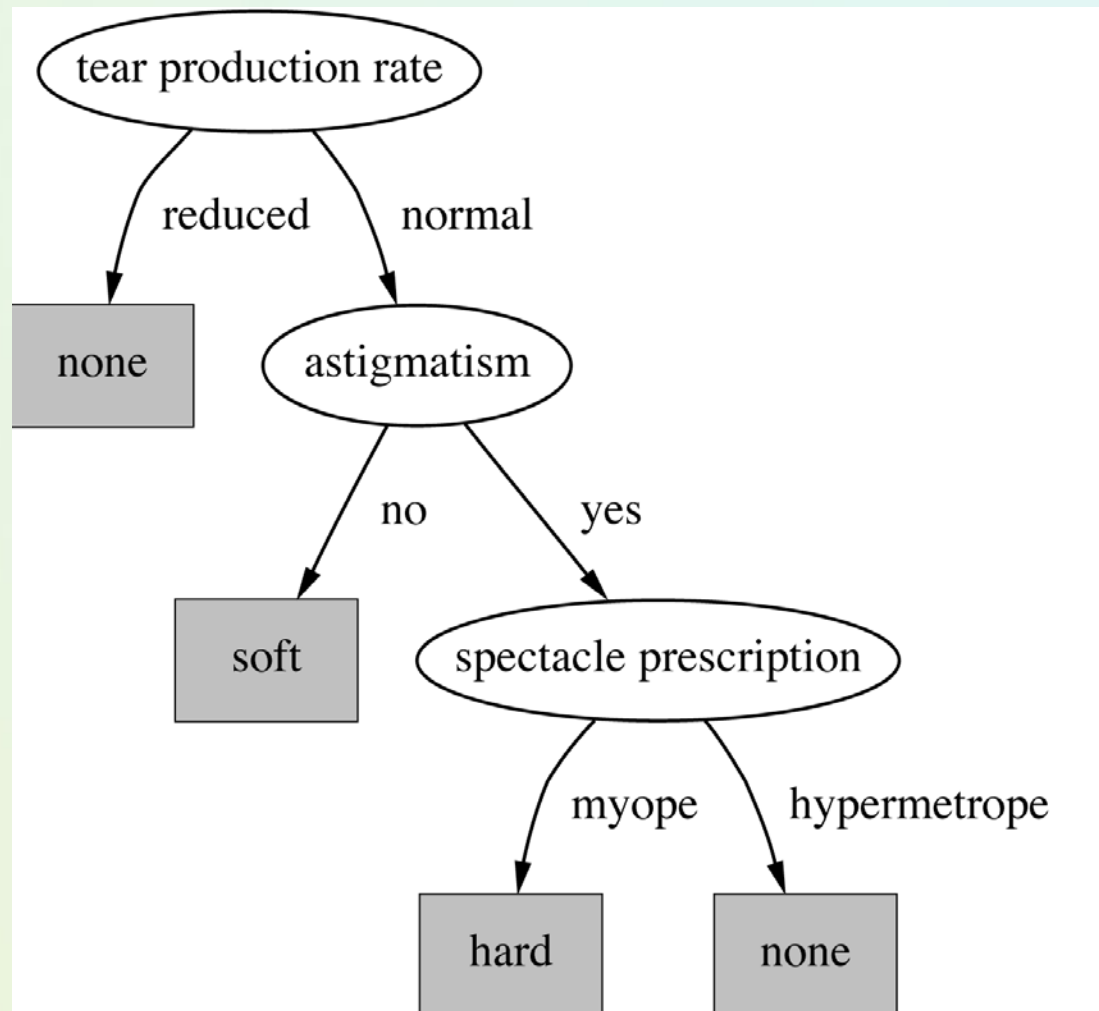
Poszukiwanie wiedzy klasyfikacyjnej (2)

- **Uczenie nadzorowane – odkrywanie wiedzy klasyfikacyjnej**
 - Przykłady uczące opisane etykietą wskazująca klasę decyzyjną; Na zbiorze uczącym poszukuje się reprezentacji wiedzy klasyfikacyjnej
- **Perspektywy odkrywania wiedzy**
 - **Predykcja** – przewidywanie przydziału nowych obiektów do klas / reprezentacja wiedzy wykorzystywana jako tzw. **klasyfikator**
 - **Opis klasyfikacji obiektów** – wyszukiwanie wzorców charakteryzujących właściwości danych i prezentacja ich użytkownikowi w zrozumiałej formie

Reprezentacja danych – tablica danych

- Tabela ($U, A \cup \{d\}$)
- Przykład tzw. contact lenses / dobór szkieł kontaktowych:
- Atrybuty:
 - age {young, pre-presbyopic, presbyopic}
 - spectacle-prescrip {myope, hypermetrope}
 - astigmatism {no, yes}
 - tear-prod-rate {reduced, normal}
- Decyzja contact-lenses {soft, hard, none}

age	specpres	astig	tearprod	contlen
young	myope	no	reduced	none
young	myope	no	normal	soft
young	myope	yes	reduced	none
young	myope	yes	normal	hard
young	hypermetrope	no	reduced	none
young	hypermetrope	no	normal	soft
young	hypermetrope	yes	reduced	none
young	hypermetrope	yes	normal	hard
pre-presbyopic	myope	no	reduced	none
pre-presbyopic	myope	no	normal	soft
pre-presbyopic	myope	yes	reduced	none
pre-presbyopic	myope	yes	normal	hard
pre-presbyopic	hypermetrope	no	reduced	none
pre-presbyopic	hypermetrope	no	normal	soft
pre-presbyopic	hypermetrope	yes	reduced	none
pre-presbyopic	hypermetrope	yes	normal	none
presbyopic	myope	no	reduced	none
presbyopic	myope	no	normal	none
presbyopic	myope	yes	reduced	none
presbyopic	myope	yes	normal	hard
presbyopic	hypermetrope	no	reduced	none
presbyopic	hypermetrope	no	normal	soft
presbyopic	hypermetrope	yes	reduced	none
presbyopic	hypermetrope	yes	normal	none



Decision tree for the contact lens data
Drzewo decyzyjne – drzewo klasyfikacyjne

Weka – software for data mining



- Waikato Environment for Knowledge Analysis (WEKA); developed by the Department of Computer Science, University of Waikato, New Zealand
- Data mining / Machine learning software written in Java (distributed under the GNU Public License)
- Used for research, education, and applications

<http://www.cs.waikato.ac.nz/ml/weka/>

- Ian Witten, Eibe Frank

WEKA – analiza pliku contact lenses

The screenshot displays the Weka Knowledge Explorer interface. The top menu bar includes 'Preprocess', 'Classify', 'Cluster', 'Associate', 'Select attributes', and 'Visualize'. Below the menu are buttons for 'Open file...', 'Open URL...', 'Open DB...', 'Undo', and 'Save...'. The 'Filter' section shows 'Choose None' and an 'Apply' button. The 'Current relation' section indicates 'Relation: contact-lenses' and 'Instances: 24'. The 'Attributes' list shows five attributes: 'age', 'spectacle-prescrip', 'astigmatism', 'tear-prod-rate', and 'contact-lenses'. The 'Selected attribute' section shows 'Name: age', 'Missing: 0 (0%)', 'Distinct: 3', and 'Type: Nominal'. A table below shows the distribution of 'age' values: 'young' (8), 'pre-presbyopic' (8), and 'presbyopic' (8). The 'Visualize All' button is visible, and the resulting visualization shows three stacked bar charts, each representing a different 'contact-lenses' category. Each bar is divided into three segments: cyan (top), red (middle), and blue (bottom), representing the distribution of 'age' values within that category. The status bar at the bottom shows 'Status OK' and a 'Log' button.

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter
Choose **None** Apply

Current relation
Relation: contact-lenses
Instances: 24
Attributes: 5

Attributes

No.	Name
1	age
2	spectacle-prescrip
3	astigmatism
4	tear-prod-rate
5	contact-lenses

Selected attribute
Name: age
Missing: 0 (0%)
Distinct: 3
Type: Nominal
Unique: 0 (0%)

Label	Count
young	8
pre-presbyopic	8
presbyopic	8

Colour: contact-lenses (Nom) Visualize All

Status
OK Log x 0

Start Total Comman... zimdataming... Weka-3-4 Weka GUI Cho... Weka Knowled... PL 20:32

Poszukiwanie drzew ze zbioru przykładów

The screenshot displays the Weka Knowledge Explorer interface. The 'Classifier' tab is active, showing the 'Id3' classifier selected. The 'Test options' section includes radio buttons for 'Use training set' (selected), 'Supplied test set', 'Cross-validation' (with 'Folds' set to 10), and 'Percentage split' (with '%' set to 66). Below this, there are 'Start' and 'Stop' buttons. The 'Result list' shows a single entry: '19:34:13 - trees.Id3'. The 'Classifier output' pane contains the following text:

```
=== Classifier model (full training set) ===  
Id3  
  
tear-prod-rate = reduced: none  
tear-prod-rate = normal  
| astigmatism = no  
| | age = young: soft  
| | age = pre-presbyopic: soft  
| | age = presbyopic  
| | | spectacle-prescrip = myope: none  
| | | spectacle-prescrip = hypermetrope: soft  
| astigmatism = yes  
| | spectacle-prescrip = myope: hard  
| | spectacle-prescrip = hypermetrope  
| | | age = young: hard  
| | | age = pre-presbyopic: none  
| | | age = presbyopic: none  
  
Time taken to build model: 0.01 seconds  
  
=== Evaluation on training set ===  
=== Summary ===  
  
Correctly Classified Instances      24          100   %  
Incorrectly Classified Instances    0           0   %  
Kappa statistic                    1  
Mean absolute error                0  
Root mean squared error            0  
Relative absolute error             0   %
```

The status bar at the bottom shows 'Status OK' and a 'Log' button. The Windows taskbar at the very bottom shows the Start button and several open applications, including 'Total Comman...', 'zimdatamining...', 'Weka-3-4', 'Weka GUI Cho...', and 'Weka Knowled...', with the system clock showing 20:34.

Widok drzewa

The screenshot displays the Weka Knowledge Explorer interface. The main window shows a decision tree visualization titled "Weka Classifier Tree Visualizer: 19:35:23 - trees.j48.J48 (contact-lenses)". The tree structure is as follows:

```
graph TD; A(tear-prod-rate) -- "= reduced" --> B[none (12.0)]; A -- "= normal" --> C(astigmatism); C -- "= no" --> D[soft (6.0/1.0)]; C -- "= yes" --> E(spectacle-prescrip); E -- "= myope" --> F[hard (3.0)]; E -- "= hypermetrope" --> G[none (3.0/1.0)];
```

The interface includes a sidebar with the following sections:

- Classifier:** J48 -U -M 2
- Test options:** Use training set (selected), Supplied test set, Cross-validation, Percentage split.
- (Nom) contact-lenses:** Start button.
- Result list (right-click for options):** 19:34:13 - trees.Id3, 19:35:23 - trees.j48.J48 (selected).

The status bar at the bottom shows "Status OK" and a "Log" button. The Windows taskbar at the very bottom shows the Start button and several open applications: Total Commander 6.0..., zimdataming.ppt, Weka-3-4, and 3 java processes.

Indukcja reguł – algorytm PRISM

The screenshot displays the Weka Knowledge Explorer interface. The 'Classifier' tab is active, and the 'Prism' classifier is selected. The 'Test options' section shows 'Use training set' selected, with 'Cross-validation' set to 10 folds and 'Percentage split' at 66%. The 'Classifier output' pane displays the following Prism rules:

```
Prism rules
-----
If astigmatism = no
  and tear-prod-rate = normal
  and spectacle-prescrip = hypermetrope then soft
If astigmatism = no
  and tear-prod-rate = normal
  and age = young then soft
If age = pre-presbyopic
  and astigmatism = no
  and tear-prod-rate = normal then soft
If astigmatism = yes
  and tear-prod-rate = normal
  and spectacle-prescrip = myope then hard
If age = young
  and astigmatism = yes
  and tear-prod-rate = normal then hard
If tear-prod-rate = reduced then none
If age = presbyopic
  and tear-prod-rate = normal
  and spectacle-prescrip = myope
  and astigmatism = no then none
If spectacle-prescrip = hypermetrope
  and astigmatism = yes
  and age = pre-presbyopic then none
If age = presbyopic
  and spectacle-prescrip = hypermetrope
  and astigmatism = yes then none

Time taken to build model: 0.01 seconds
```

The 'Result list' shows the following entries:

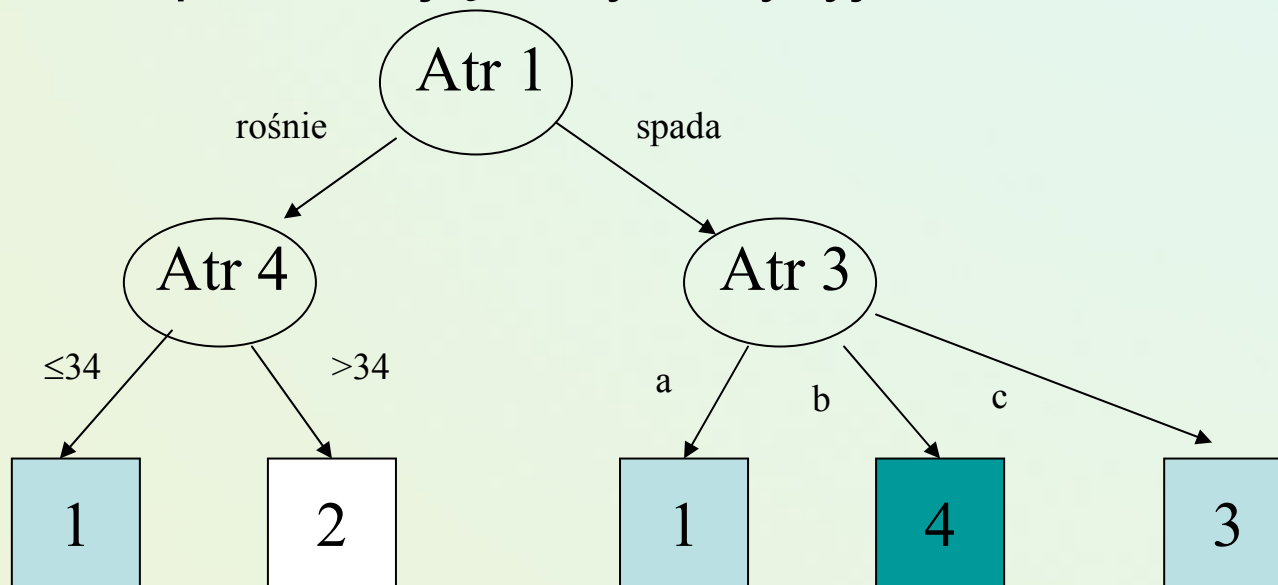
- 19:34:13 - trees.Id3
- 19:35:23 - trees.j48.J48
- 19:37:12 - rules Prism

The status bar at the bottom indicates 'OK' and shows a 'Log' button. The Windows taskbar at the very bottom shows the Start button and several open applications, including 'Total Command...', 'zimdataminging.ppt', 'Weka-3-4', 'Weka GUI Choo...', and 'Weka Knowledg...'. The system clock shows 20:37.

Co to jest drzewo decyzyjne?

Jest to struktura grafu skierowanego z góry na dół:

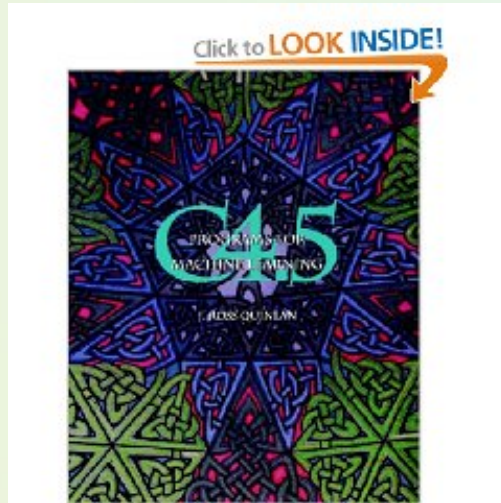
- Węzły reprezentują pytanie o wartości cech
- Z węzłów wychodzą gałęzie które reprezentują wynik pytania
- Liście reprezentują klasy decyzyjne



Metody indukcji drzew decyzyjnych

- Podejście „top-down decision tree generation” obejmuje dwa etapy:
 - **Konstrukcja drzewa**
 - Na początku wszystkie przykłady w węźle.
 - Rekurencyjnie dziel przykłady w oparciu o wybrane testy na wartościach atrybutu (kryterium wyboru).
 - Upraszczenie drzewa - „Tree pruning”
 - Usuwanie poddrzew, które mogą prowadzić do błędnych decyzji podczas klasyfikacji przypadków testowych („noisy data”, przespecjalizowane opisy ...)
- Przykłady algorytmów: ID3, C4.5, CART,...

J.Ross Quinlan – twórca alg. indukcji drzew

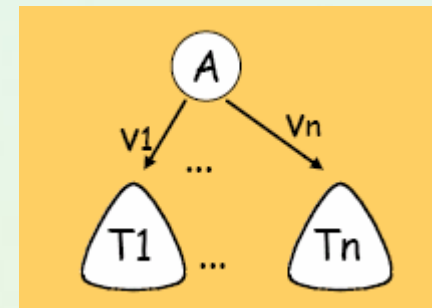


Ross Quinlan completed a PhD in Computer Science at the University of Washington in 1968. He has developed several algorithms used in machine learning and data mining such as ID3, C4.5, FOIL, and more recent commercial systems such as See5 and Cubist. He has held permanent appointments at the University of Sydney, University of Technology Sydney, Rand Corporation, and visiting appointments at Carnegie-Mellon University, MIT, GTE, and Stanford University. He currently heads a small data mining tools company and is an Adjunct Professor at the University of New South Wales. He is a Fellow of the American Association for Artificial Intelligence and the Australian Computer Society.

- Więcej – spójrz <http://www.rulequest.com/Personal/>
- Także, http://en.wikipedia.org/wiki/Ross_Quinlan

Basic TDIDT algorithm (simplified Quinlan's ID3)

- At start, all training examples S are at the root.
- **If** all examples from S belong to the same class K_j
then label the root with K_j
else
 - select the „best” attribute A
 - divide S into S_1, \dots, S_n according to values v_1, \dots, v_n of attribute A
 - Recursively build subtrees T_1, \dots, T_n for S_1, \dots, S_n

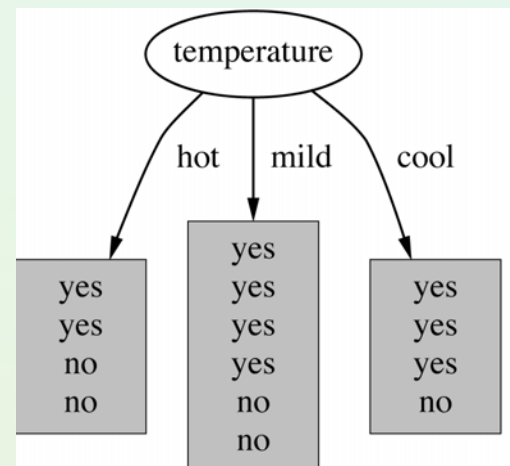
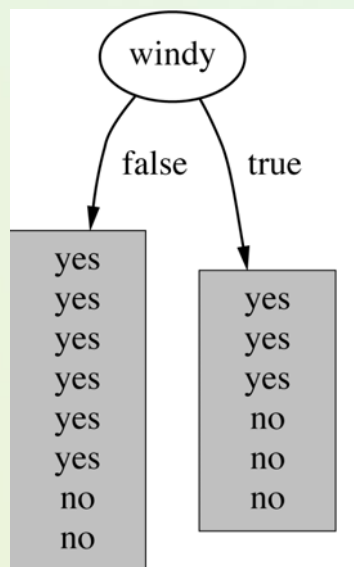
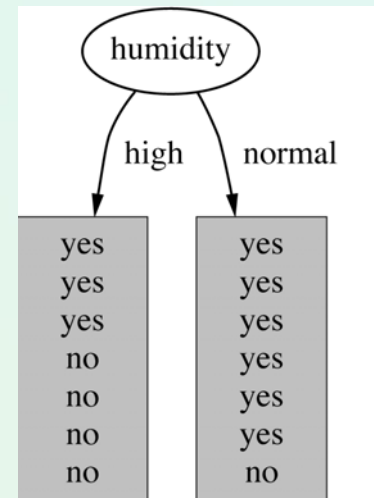
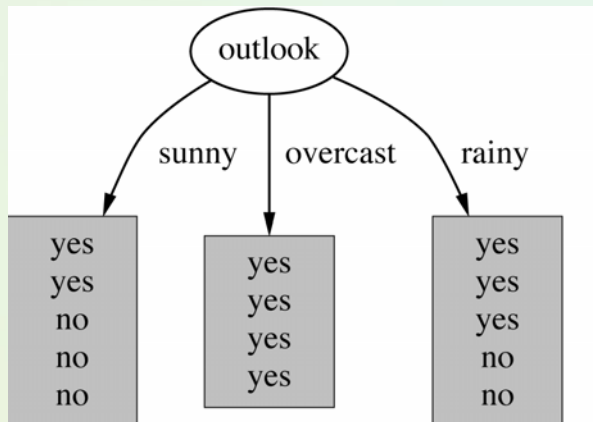


Przykład budowy – Quinlan example „golf”

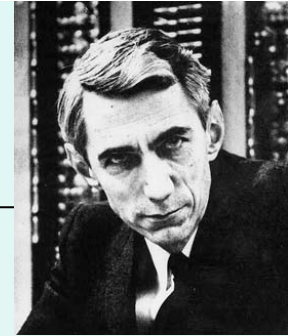
Outlook	Temperature	Humidity	Windy	Play?
sunny	hot	high	false	No
sunny	hot	high	true	No
overcast	hot	high	false	Yes
rain	mild	high	false	Yes
rain	cool	normal	false	Yes
rain	cool	normal	true	No
overcast	cool	normal	true	Yes
sunny	mild	high	false	No
sunny	cool	normal	false	Yes
rain	mild	normal	false	Yes
sunny	mild	normal	true	Yes
overcast	mild	high	true	Yes
overcast	hot	normal	false	Yes
rain	mild	high	true	No

*Uproszczona
Tabela danych*

Który atrybut należy wybrać?



Entropia informacji



- Entropy information (wprowadził C. Shannon)
 - Given a probability distribution, the info required to predict an event is the distribution's *entropy*
 - Entropy gives the information required in bits (this can involve fractions of bits!)
 - The amount of information, needed to decide if an arbitrary example in S belongs to class K_j (p_j - prob. it belongs to K_j).

- Basic formula for computing the entropy for examples in S :

$$\text{entropy}(S) = -p_1 \log p_1 - p_2 \log p_2 \dots - p_n \log p_n$$

- A conditional entropy for splitting examples S into subsets S_i by using an attribute A :

$$\text{entropy}(S | A) = \sum_{i=1}^m \frac{|S_i|}{|S|} \cdot \text{entropy}(S_i)$$

- Choose the attribute with the maximal info gain:

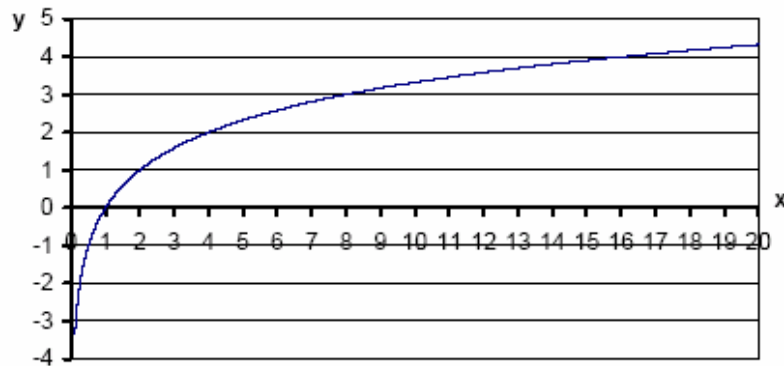
$$\text{entropy}(S) - \text{entropy}(S | A)$$

Funkcja logarytmiczna

$$y = \log_a x$$

a – podstawa logarytmu $x = a^y$

Rozważmy funkcję logarytmiczną dla $a = 2$ (tj. $y = \log_2 x$)



x	1/8	1/4	1/2	1	2	4	8
y	-3	-2	-1	0	1	2	3

- Co zrobić, gdy $p=0$
- Jak wygląda wykres entropii dla klasyfikacji binarnej

Entropia dla przykładu golf

Nie oceniamy podziału atrybutem, tylko rozkład wartości klas decyzyjnych

Dwie klasy : *yes* and *no*

Z 14 przykładów 9 etykietowanych jako *yes*, reszta jako *no*

$$p_{yes} = -\left(\frac{9}{14}\right) \log_2 \left(\frac{9}{14}\right) = 0.41$$

$$p_{no} = -\left(\frac{5}{14}\right) \log_2 \left(\frac{5}{14}\right) = 0.53$$

$$E(S) = p_{yes} + p_{no} = 0.94$$

Outlook	Temp.	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes

Outlook	Temp.	Humidity	Windy	play
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Przykład oceny atrybutu "Outlook"

- "Outlook" = "Sunny":

$$\text{info}([2,3]) = \text{entropy}(2/5,3/5) = -2/5 \log(2/5) - 3/5 \log(3/5) = 0.971$$

- "Outlook" = "Overcast":

$$\text{info}([4,0]) = \text{entropy}(1,0) = -1 \log(1) - 0 \log(0) = 0$$

Note: $\log(0)$ is not defined, but we evaluate $0 \cdot \log(0)$ as zero

- "Outlook" = "Rainy":

$$\text{info}([3,2]) = \text{entropy}(3/5,2/5) = -3/5 \log(3/5) - 2/5 \log(2/5) = 0.971$$

- Entropia warunkowa dla podziału wartościami atrybutu

$$\begin{aligned} \text{info}([3,2],[4,0],[3,2]) &= (5/14) \times 0.971 + (4/14) \times 0 + (5/14) \times 0.971 \\ &= 0.693 \end{aligned}$$

Obliczanie zysku informacyjnego miary entropii

- Information gain:

(information before split) – (information after split)

$$\text{gain("Outlook")} = \text{info}([9,5]) - \text{info}([2,3],[4,0],[3,2]) = 0.940 - 0.693 \\ = 0.247$$

- Ostateczne wartości zysku

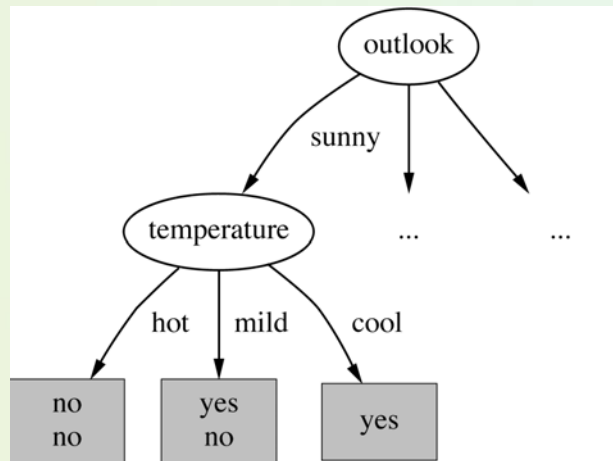
$$\text{gain("Outlook")} = 0.247$$

$$\text{gain("Temperature")} = 0.029$$

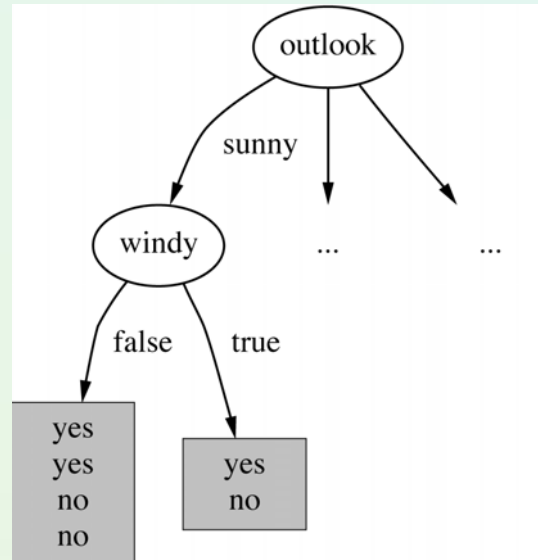
$$\text{gain("Humidity")} = 0.152$$

$$\text{gain("Windy")} = 0.048$$

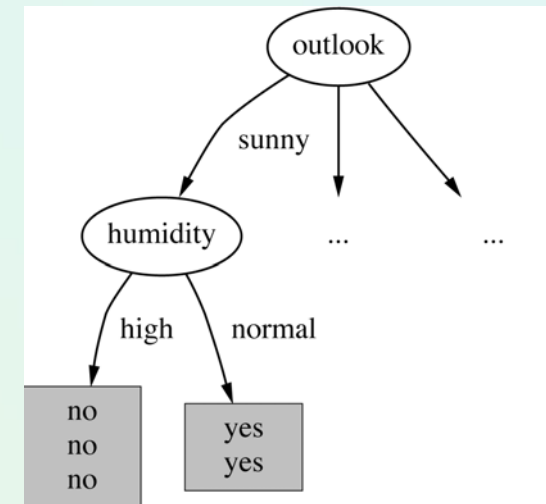
Dalsze obliczenia



$$\text{gain}(\text{"Temperature"}) = 0.571$$

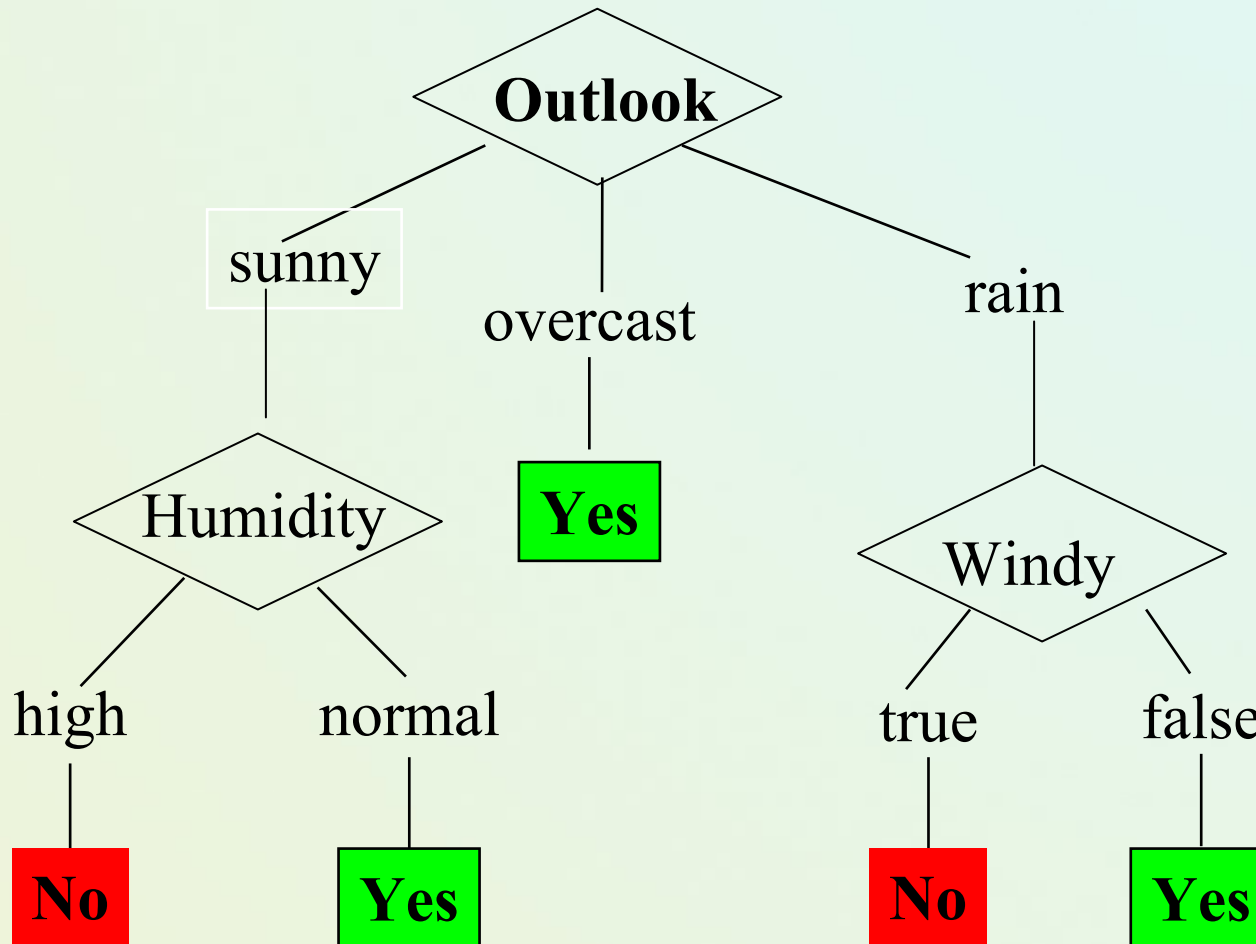


$$\text{gain}(\text{"Windy"}) = 0.020$$



$$\text{gain}(\text{"Humidity"}) = 0.971$$

Ostateczne drzewo

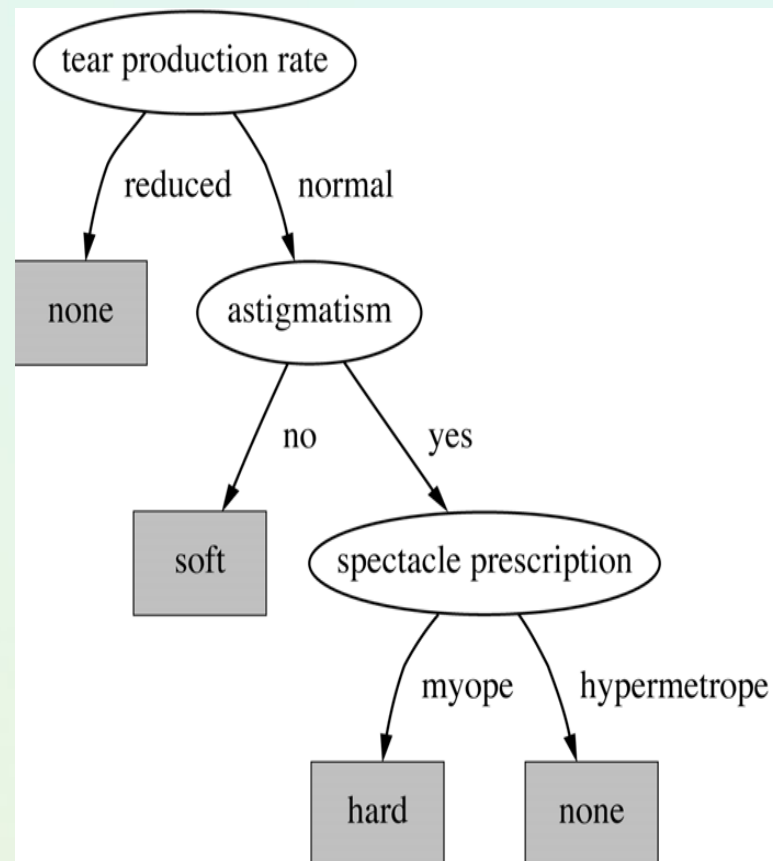


Podstawowe pojęcia w indukcji DT

- Kryterium podziału: *Goodness function*
 - Stosowana, aby wybrać atrybut do tworzenia węzła
 - Różne funkcje są używane:
 - information gain (*entropia*)
 - Gini index
- Tworzenie gałęzi drzewa:
 - Ustalanie gałęzi, do której przydzielamy podzbiór.
 - binary versus k -ary splitting.
- Decyzja kiedy zatrzymać rozbudowę drzewa: impurity measure.
- Tworzenie liści: przypisanie etykiety klasy większościowej.

Wykorzystanie drzewa

- Bezpośrednio:
 - sprawdzaj wartości atrybutu nowego przykładu zaczynając od korzenia do liści
- Pośrednio:
 - zamień strukturę drzewa na zbiór reguł decyzyjnych (upraszczając nadmiarowe warunki)
 - reguły uważa się za czytelniejszą reprezentację



Ograniczenia w uczeniu się drzew decyzyjnych

Pytania i problemy, np.:

Kiedy należy zaprzestać rozbudowywać drzewa?

aby zapobiec przespecjalizowaniu opisu

duże drzewa są trudne do analizy i zrozumienia

Jak uwzględniać atrybuty ilościowe?

Jak uwzględniać atrybuty ze zbyt dużą liczbą wartości w stosunku do dziedzin pozostałych atrybutów?

Jak uwzględniać atrybuty z nieznanymi wartościami?

Jak uwzględniać dane "zaszumione"?

Binary Tree – budowa drzew binarnych

- Drzewa binarne mogą być skuteczniejsze w klasyfikacji nowych faktów
- Podział binarny w węźle drzewa:
 - Atrybuty liczbowe A , reprezentacja w postaci $\text{value}(A) < x$ gdzie x jest wartością z dziedziny A .
 - Atrybuty nieliczbowe A , warunek w postaci $\text{value}(A) \in X$ gdzie $X \subset \text{domain}(A)$

Drzewo binarne (Quinlan's C4.5 output)

Pruned decision tree:

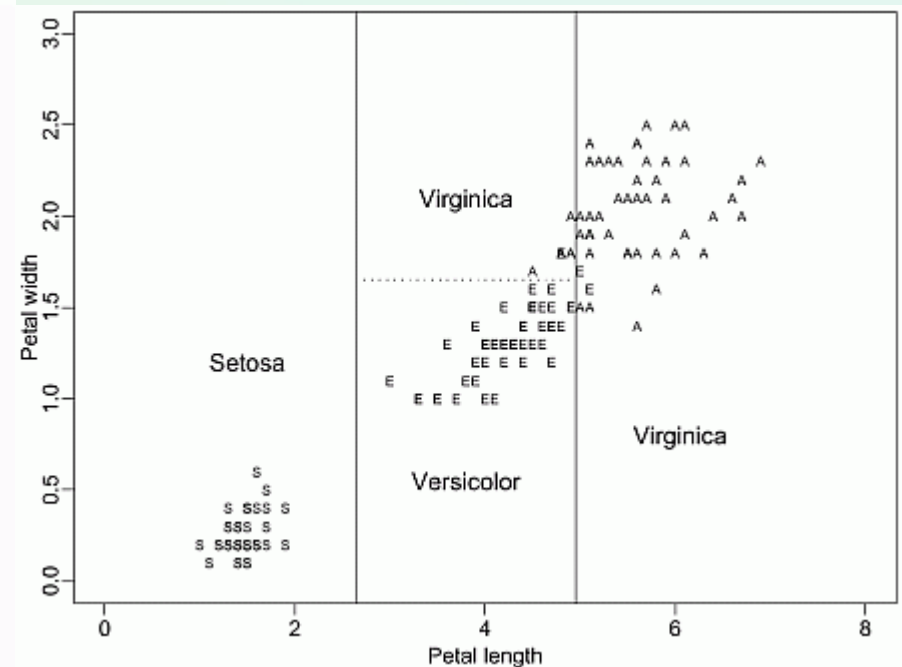
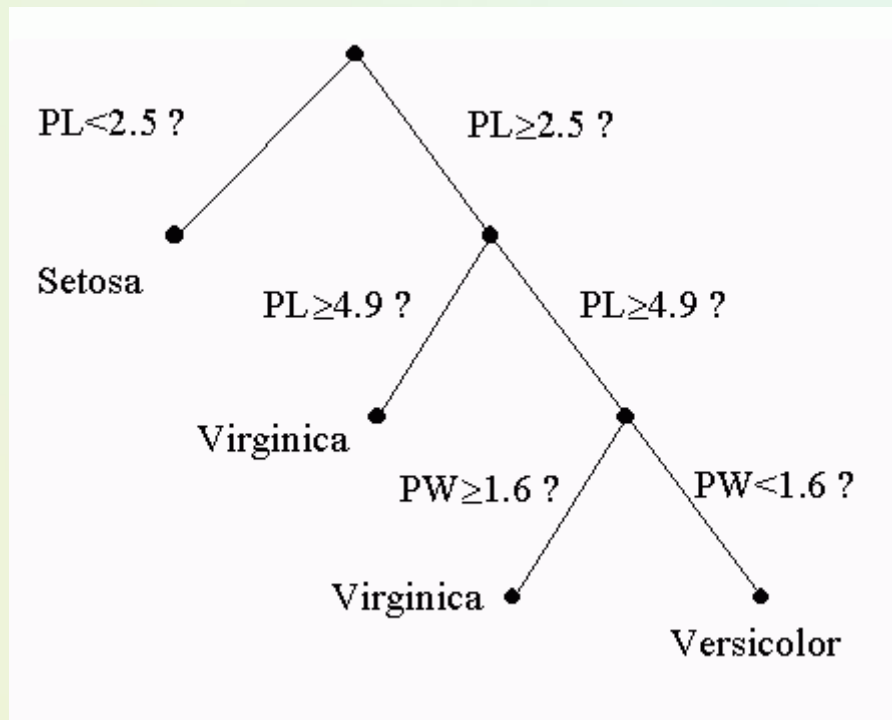
```
A9 - t:
  A15 > 228 : + (106.0/3.8)
  A15 <= 228 :
    A14 <= 102 :
      A4 in {l,t}: + (0.0)
      A4 - u:
        A6 in {c,d,cc,i,k,m,q,w,x,e,aa}: + (46.4/3.1)
        A6 in {j,ff}: - (2.0/1.0)
        A6 - r: + (0.0)
      A4 - y:
        A6 in {c,i,aa,ff}: - (7.0/3.4)
        A6 in {d,j,w,x}: + (4.0/1.2)
        A6 in {cc,k,m,r,q,e}: + (0.0)
    A14 > 102 :
      A6 in {j,r}: + (0.0)
      A6 in {c,d,k,m,e,aa,ff}:
        A14 <= 132 : - (4.1/1.2)
        A14 > 132 :
          A3 <= 1.625 :
            A14 <= 292 : - (13.0/1.3)
            A14 > 292 :
              A13 - g: + (2.0/1.0)
              A13 - s: - (6.0/2.3)
              A13 - p: - (0.0)
          A3 > 1.625 :
            A6 in {k,m}: + (5.0/1.2)
            A6 - ff: + (0.0)
            A6 in {c,d,e,aa}:
              A2 <= 32.08 : + (9.5/4.1)
              A2 > 32.08 : - (8.0/3.5)
            A6 in {cc,i,q,w,x}:
              A8 <= 10.75 : + (36.0/9.3)
              A8 > 10.75 : - (2.0/1.0)
  A9 - f:
    A4 in {u,y}: - (237.0/17.3)
    A4 - l: + (2.0/1.0)
    A4 - t: - (0.0)
```

- Crx (Credit Data) UCI ML Repository

Interpretacja graficzna dla atr liczbowych

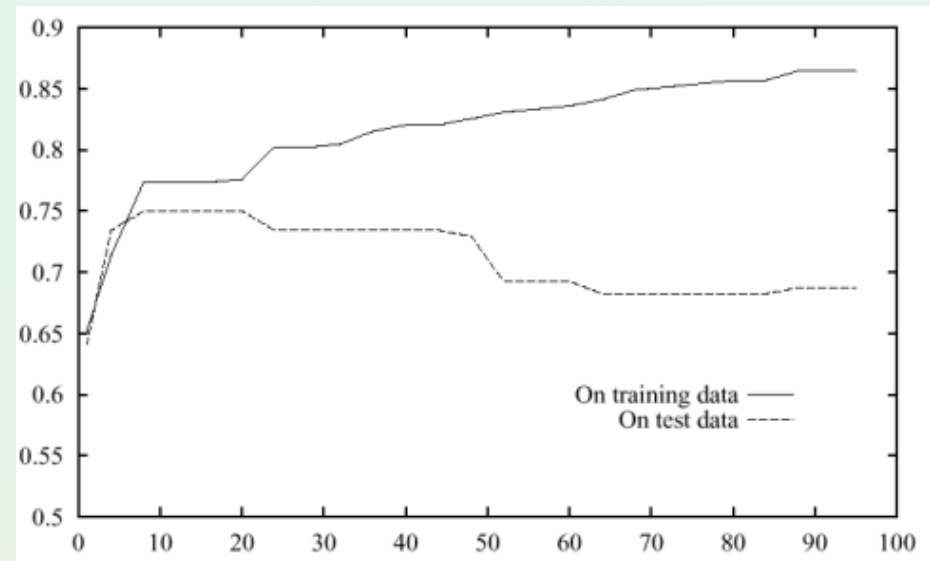
Hierarchiczny podział na hiper-prostopadłościany

Przykład: Iris flowers data, with 4 features; displayed in 2-D



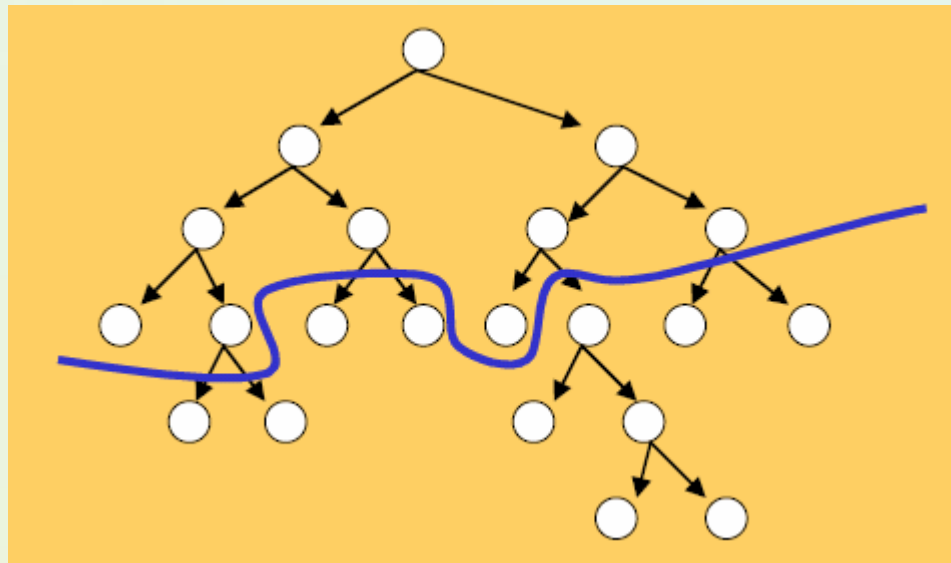
Overfitting the Data – nadmierne dopasowanie do danych uczących

- Podstawowy algorytm ID3 → Rozbuduj gałąź drzewa do pełnego rozróżnienia przykładów
 - Sensowe na spójnych przykładów i celów dokładnego opisu
- Rzeczywiste dane (nieśpójne, szum informacyjny) oraz cel klasyfikowania przykładów
 - Drzewa mają tendencje do przeuczenia / nadmiernego dopasowania do specyficznych przykładów *overfit* the learning examples
 - Occam razor – zasada brzytwy Occama (z konkurujących drzew wybierz prostsze; ma lepsze własności generalizacyjne)



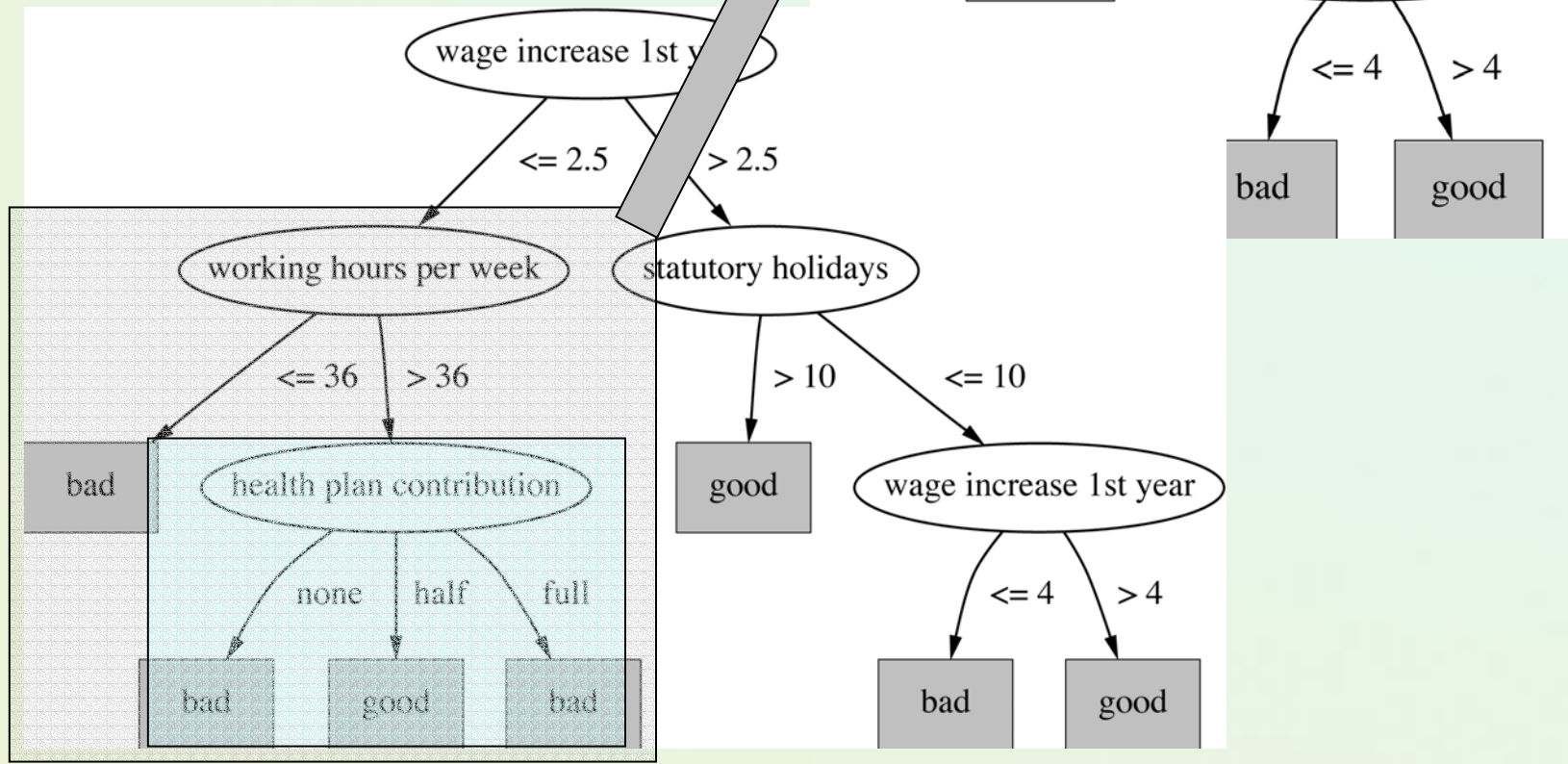
Tree pruning – upraszczanie drzewa

- Tree pruning – mechanizm „walki” z przeuczeniem
- Po uproszczeniu struktury drzewa może wzrosnąć trafność na przykładach testowych!



Przykład redukcji

- *Tzw. post-pruning*
- Usuwać podzrewa i oceniaj wpływ na estymatę błędu / poprawności decyzji klasyfikacyjnych



Przykłady zastosowań w medycynie

- Wiele przykładów analizy podejmowania decyzji o diagnozowaniu chorób, także terapii, oraz farmacja i budowa związków - leków:
- Przykładowe omówienia:
 - I.Kononenko, I.Bratko, M.Kukar: Application of Machine Learning to Medical Diagnosis. w: Michalski R.S., Bratko I, Kubat M. (red.), Machine learning and data mining, John Wiley & Sons, 1998, s. 389-408.
 - Langley, P., Simon, H. A., Fielded applications of machine learning, w: Michalski R.S., Bratko I, Kubat M. (red.), Machine learning and data mining, John Wiley & Sons, 1998 , s. 113-129.
- Spójrz także na dodatkowe slajdy (włączmy tradycyjny rzutnik 😊...)

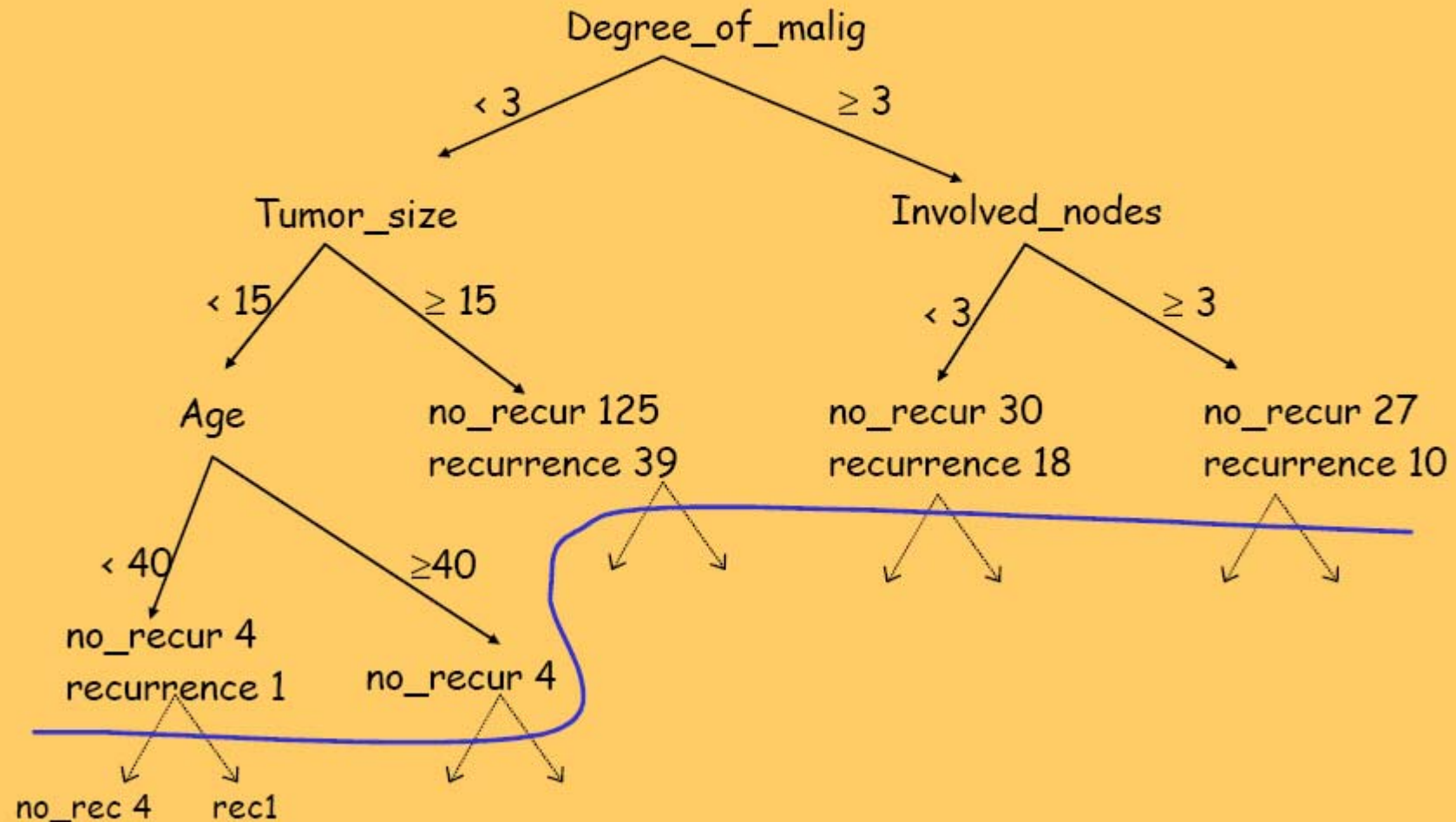
Przykład drzewa decyzyjnego

Medicine - Predicting C-Section Risk

- Learned from Medical Records of 1000 Women
- Negative Examples are Cesarean Sections
 - Prior distribution: [833+, 167-] 0.83+, 0.17-
 - *Fetal-Presentation* = 1: [822+, 167-] 0.88+, 0.12-
 - *Previous-C-Section* = 0: [767+, 81-] 0.90+, 0.10-
 - *Primiparous* = 0: [399+, 13-] 0.97+, 0.03-
 - *Primiparous* = 1: [368+, 68-] 0.84+, 0.16-
 - *Fetal-Distress* = 0: [334+, 47-] 0.88+, 0.12-
 - *Birth-Weight* < 3349 0.95+, 0.05-
 - *Birth-Weight* ≥ 3347 0.78+, 0.22-
 - *Fetal-Distress* = 1: [34+, 21-] 0.62+, 0.38-
 - *Previous-C-Section* = 1: [55+, 35-] 0.61+, 0.39-
 - *Fetal-Presentation* = 2: [3+, 29-] 0.11+, 0.89-
 - *Fetal-Presentation* = 3: [8+, 22-] 0.27+, 0.73-

A real life example of tree pruning [U.Ljubljana]

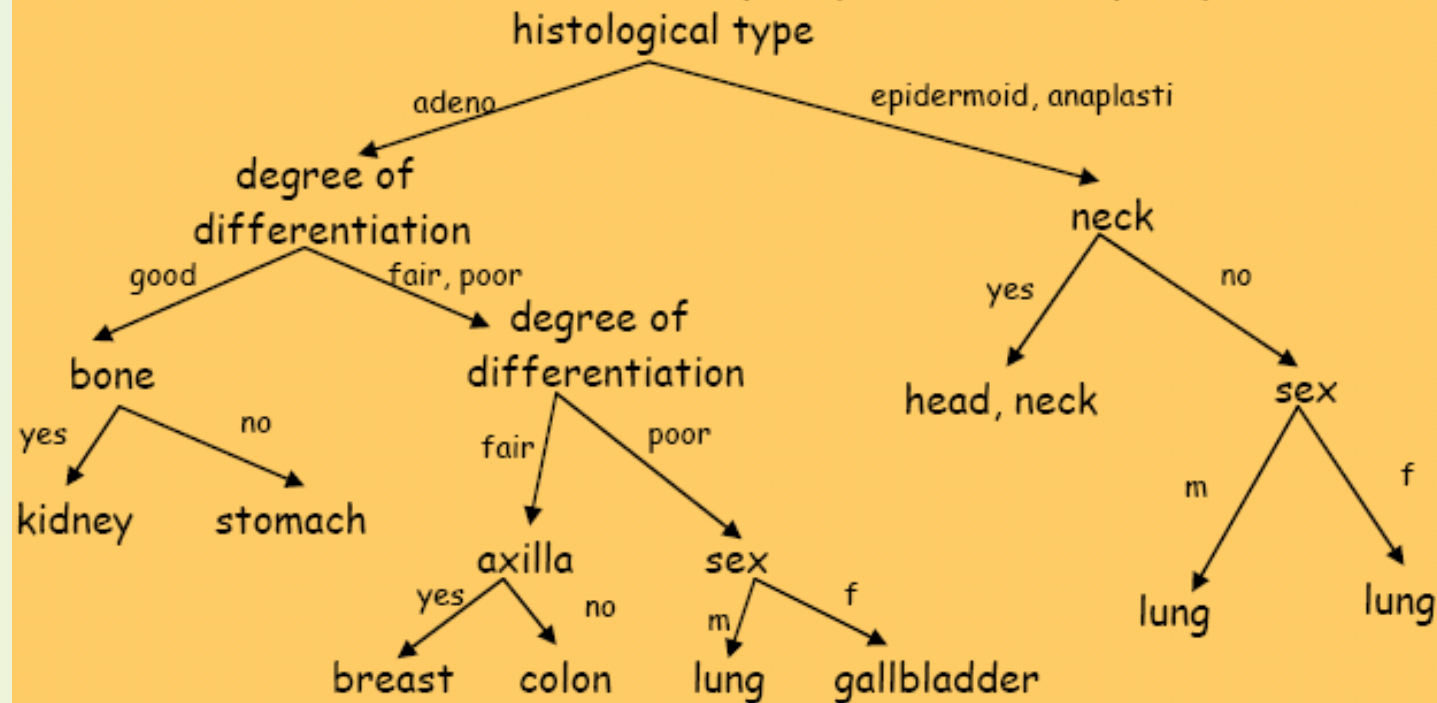
Prediction of breast cancer recurrence: Tree pruning



Another example from U.Ljubljana

Location of primary tumor

- 339 examples
- 228 for learning, 111 for testing
- induce decision tree accuracy: unpruned: 41%, postpruned: 45%



C4.5 Quinlan + PP interfejs

- Opiera się na oryginalnym kodzie J.R. Quinlana

The screenshot displays the C4.5 software interface with several windows open. The main window, titled "C4.5 GOLF (4 attributes, 14 training cases, 11 test cases)", shows a table comparing the performance of a tree before and after pruning.

Before pruning				After pruning			
Tree	Size	Errors	Errors (test)	Size	Errors	Errors (test)	Estimate
1	8	0 (0.0%)	1 (9.1%)	8	0 (0.0%)	1 (9.1%)	38.5%

Two confusion matrix windows are also visible:

- Confusion matrix (training set):**

Org. \ C4.5	Play	Don't Play
Play	9	
Don't Play		5
- Confusion matrix (test set):**

Org. \ C4.5	Play	Don't Play
Play	7	
Don't Play	1	3

At the bottom, two windows show the decision tree structure:

- Unpruned tree:** A tree with root node "outlook = overcast" leading to "Play". Other branches include "outlook = sunny" (with a humidity split) and "outlook = rain" (with a windy split).
- Pruned tree:** A tree with root node "outlook = overcast" leading to "Play". Other branches include "outlook = sunny" (with a humidity split) and "outlook = rain" (with a windy split).

The Windows taskbar at the bottom shows the Start button, Total Commander 5.5..., C4.5, and system tray icons including the date and time (16:45).

CART from Salford Systems

Salford Systems - Windows Internet Explorer

E:\zdjecia\materialy\Breiman\products-cart.html

Live Search

Plik Edycja Widok Ulubione Narzędzia Pomoc

Salford Systems

Strona Narzędzia

SALFORD SYSTEMS Home - [Sitemap](#) - [More Info](#) - [Contact us](#) - [Help](#)

Products Services Resources News & Events Support Company Salford Conference Success Stories

Latest news

- [CART 6.0 ProEX-- New For 2008](#)
- [Salford Tools Win 2007 DMA Challenge, PAKDD Competition](#)
- [User Group Case Study Presentations Now Available](#)
- [Announcing Salford Systems Success Stories](#)
- [TreeNet Wins Major Competition](#)

Events


- [Hands-On Training in San Diego: December 2008](#)
- [Salford Systems Announces 2009 Conference](#)
- [Previous Conference Proceedings Now Available](#)

Highlights


- [Read about how Fleet Financial Group used CART® to improve their customer service](#)


Home Salford Systems

Leads the business intelligence and data mining industries by converting significant new scientific discoveries into widely accessible, high performance software solutions.

 Check out how our award-winning products such as **CART®**, **MARS®**, **TreeNet®**, and **RandomForests™** can help you predict the future of your business NOW!

Speak to our team of in-house **consultants** to see how they can help you build practical data mining solutions.



 Need to determine your ROI on your data mining project quickly and cost-effectively? Consider our **Rapid Response Data Mining Center**.

Winner of the 2003 predictive modeling and data mining tournament at **Teradata**

Analiza dużych bazach danych

- Klasyfikacja rozważana głównie przez:
 - statystyków
 - badaczy z AI – „machine learning”
- Data mining wymaga rozważenia problemów związanych z przetwarzaniem dużych baz danych
 - dotychczasowe zastosowania na małych rozmiarach danych, algorytmy wykorzystują PAO
- Data mining dostarcza rozszerzeń w zakresie:
 - Skalowaności,
 - Przetwarzania równoległego i rozproszonego,
 - Danych o złożonej strukturze,
 -

Drzewa decyzyjne a data mining ?

- Wybór drzew decyzyjnych:
 - Względnie szybkie algorytmy indukcji drzew, z możliwością modyfikacji dla dużych zbiorów danych
 - Możliwość transformacji w reguły klasyfikacyjne
 - Jawna reprezentacja wiedzy w postaci symbolicznej wygodnej do interpretacji przez człowieka
 - Łatwość realizacji klasyfikowania nowych przypadków
 - Osiągają wysoką trafność klasyfikowania, porównywalną z innymi metodami
 - Odporność na niedoskonałe dane
- Zastosowania do problemów z bardzo dużą liczbą przykładów opisanych setkami atrybutów (akceptowalne czasy odpowiedzi)

Indukcja reguł decyzyjnych

- Podstawowa idea - reguły poszukuje się **bezpośrednio z danych**
 - potencjalnie większa zrozumiałość wiedzy
 - ale więcej różnych podejść
 - opisy rzeczywistych danych minimalnym zbiorem reguł dyskryminujących/ klasyfikujących
 - poszukiwanie bardziej wyczerpujących zbiorów reguł o dobrych własnościach interpretacyjnych
 - więcej parametrów do sterowania w metodach regułowych

Indukcja reguł metodą generowania kolejnych pokryć

Sequential covering (X_j klasa; A atrybuty; E przykłady, τ próg akceptacji);

begin

$R := \emptyset;$ {zbiór poszukiwanych reguł}

$r := \text{learn-one-rule}(\text{klasa } X_j; \text{ atrybuty } A; \text{ przykłady } E)$

while $\text{evaluate}(r, E) > \tau$ **do**

begin

$R := R \cup r;$

$E := E \setminus [R];$ {usuń przykłady pozytywne pokryte przez R }

$r := \text{learn-one-rule}(\text{klasa } X_j; \text{ atrybuty } A; \text{ przykłady } E);$

end;

return R

end.

- Funkcja *learn-one-rule* dla danego zbioru przykładów znajduje jedną regułę pokrywającą możliwie jak najwięcej przykładów pozytywnych i jak najmniej negatywnych.



Ryszard Michalski

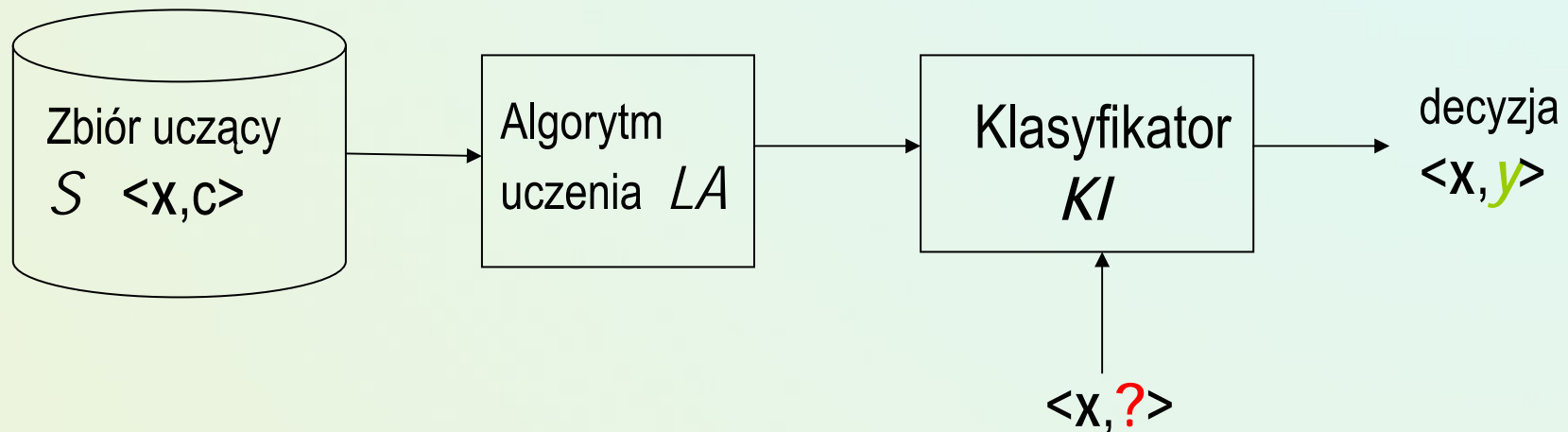
MODLEM – algorytm indukcji reguł

- MODLEM [Stefanowski 98] generuje minimalny zbiór reguł pokrywających zbiór przykładów uczących
- Właściwości – możliwość wykorzystanie teorii zbiorów przybliżonych do analizy sprzecznych przykładów; bezpośrednie przetwarzanie atrybutów liczbowych (bez wstępnej dyskretyzacji) oraz nieznanymi wartości atrybutów.
- Strategia wykorzystywania reguł → LERS [Grzymała 94]

obj.	a_1	a_2	a_3	a_4	D	
x_1	m	2.0	1	a	C1	<i>if</i> ($a_1 = m$) <i>and</i> ($a_2 \leq 2.6$) <i>then</i> ($D = C1$) { x_1, x_3, x_7 }
x_2	f	2.5	1	b	C2	<i>if</i> ($a_2 \in [1.45, 2.4]$) <i>and</i> ($a_3 \leq 2$) <i>then</i> ($D = C1$)
x_3	m	1.5	3	c	C1	{ x_1, x_4, x_7 }
x_4	f	2.3	2	c	C1	<i>if</i> ($a_2 \geq 2.4$) <i>then</i> ($D = C2$) { x_2, x_6 }
x_5	f	1.4	2	a	C2	<i>if</i> ($a_1 = f$) <i>and</i> ($a_2 \leq 2.15$) <i>then</i> ($D = C2$) { x_5, x_8 }
x_6	m	3.2	2	c	C2	
x_7	m	1.9	2	b	C1	
x_8	f	2.0	3	a	C2	

Predykcja nowych faktów - klasyfikatory

- Poszukiwanie reprezentacji wiedzy o przydziale obiektów do klas na podstawie opisu obiektów za pomocą wartości atrybutów (zbiór uczący).
- **Predykcja klasyfikacji** nowych obiektów (zbiór testowy)



Przykłady $S = \{ \langle \mathbf{x}_1, c_1 \rangle, \langle \mathbf{x}_2, c_2 \rangle, \dots, \langle \mathbf{x}_n, c_n \rangle \}$
 $\mathbf{x}_i = \langle x_{i1}, x_{i2}, \dots, x_{im} \rangle$ opisywane przez m atrybutów
Atrybuty różnego typu
 c_i – etykieta jednej z klas $\{C_1, \dots, C_K\}$

Miara oceny, np:
trafność klasyfikowania

$$\eta = \frac{N_c}{N_t}$$

Eksperymentalna ocena
→ *Cross validation*

Trafność klasyfikowania

- Użyj przykładów testowych nie wykorzystanych w fazie indukcji klasyfikatora:
 - N_t – liczba przykładów testowych
 - N_c – liczba poprawnie sklasyfikowanych przykładów testowych
- Trafność klasyfikowania (classification accuracy):

$$\eta = \frac{N_c}{N_t}$$

- Alternatywnie błąd klasyfikowania.

$$\varepsilon = \frac{N_t - N_c}{N_t}$$

Inne możliwości analizy:

- macierz pomyłek (ang. confusion matrix),
- koszty pomyłek i klasyfikacja binarna,
- miary Sensitivity i Specificity / krzywa ROC

Macierz pomyłek

- Analiza pomyłek w przydziale do różnych klas przy pomocy tzw. macierz pomyłek (ang. *confusion matrix*)
- Macierz $r \times r$, gdzie wiersze odpowiadają poprawnym klasom decyzyjnym, a kolumny decyzjom przewidywanym przez klasyfikator; na przecięciu wiersza i oraz kolumny j - liczba przykładów n_{ij} należących oryginalnie do klasy i -tej, a zaliczonej do klasy j -tej

Przykład:

	Przewidywane klasy decyzyjne		
Oryginalne klasy	K_1	K_2	K_3
K_1	50	0	0
K_2	0	48	2
K_3	0	4	46

Klasyfikacja binarna

- Niektóre zastosowania → jedna z klas posiada szczególne znaczenie, np. diagnozowanie poważnej choroby. Problem → klasyfikacja binarna.

Oryginalne klasy	Przewidywane klasy decyzyjne	
	Pozytywna	Negatywna
Pozytywna	<i>TP</i>	<i>FN</i>
Negatywna	<i>FP</i>	<i>TN</i>

- Nazewnictwo (inspirowane medycznie):
 - TP* (ang. *true positive*) – liczba poprawnie sklasyfikowanych przykładów z wybranej klasy (ang. *hit*),
 - FN* (ang. *false negative*) – liczba błędnie sklasyfikowanych przykładów z tej klasy, tj. decyzja negatywna podczas gdy przykład w rzeczywistości jest pozytywny (błąd pominięcia - z ang. *miss*),
 - TN* (ang. *true negative*) – liczba przykładów poprawnie nie przydzielonych do wybranej klasy (poprawnie odrzuconych z ang. *correct rejection*),
 - FP* (ang. *false positive*) – liczba przykładów błędnie przydzielonych do wybranej klasy, podczas gdy w rzeczywistości do niej nie należą (ang.

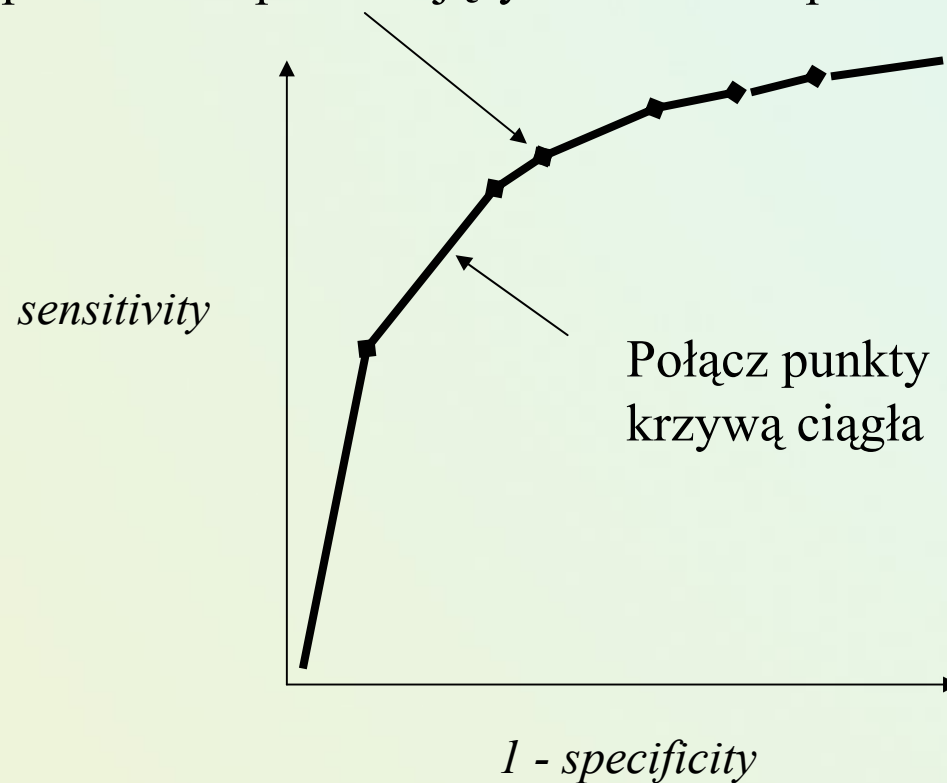
Miary stosowane w analizie klasyfikacji binarnej

- Dodatkowe miary oceny rozpoznawania wybranej klasy:
 - Wrażliwość / czułość (ang. *sensitivity*) = $TP / (TP+FN)$,
 - Specyficzność (ang. *specificity*) = $TN / (FP+TN)$.
- Inne miary:
 - *False-positive rate* = $FP / (FP+TN)$, czyli 1 – specyficzność.
- Wnikliwszą analizę działania klasyfikatorów binarnych dokonuje się w oparciu o analizę krzywej ROC, ang. *Receiver Operating Characteristic*).

Oryginalne klasy	Przewidywane klasy decyzyjne	
	Pozytywna	Negatywna
Pozytywna	<i>TP</i>	<i>FN</i>
Negatywna	<i>FP</i>	<i>TN</i>

Krzywa ROC - analiza

Algorytm może być parametryzowany, i w rezultacie otrzymuje się serie punktów odpowiadających doborowi parametrów



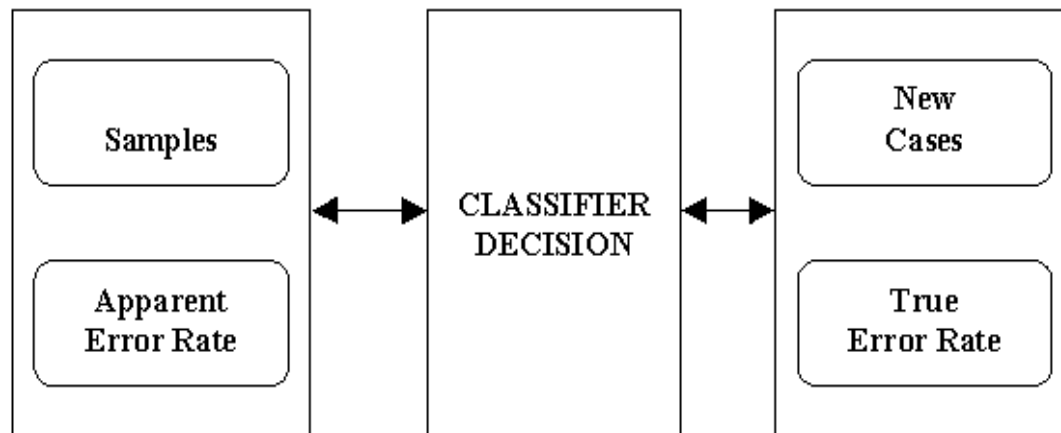
Wykres nazywany
'krzywą' ROC.

Jak szacować wiarygodnie ?

- Zależy od perspektywy użycia wiedzy:
 - Predykcja klasyfikacji albo opisowa
- Ocena na zbiorze uczącym nie jest wiarygodna jeśli rozważamy predykcję nowych faktów!
 - Nowe obserwacje najprawdopodobniej nie będą takie same jak dane uczące!
 - Choć zasada reprezentatywności próbki uczącej ...
- Problem przeuczenia (ang. overfitting)
 - Nadmierne dopasowanie do specyfiki danych uczących powiązane jest najczęściej z utratą zdolności uogólniania (ang. generalization) i predykcji nowych faktów!

Podójście empiryczne

- Zasada „Train and test”
- Gdy nie ma podziału zadanego przez nauczyciela, to wykorzystaj losowe podziały
- Nadal pytanie jak szacować wiarygodnie?

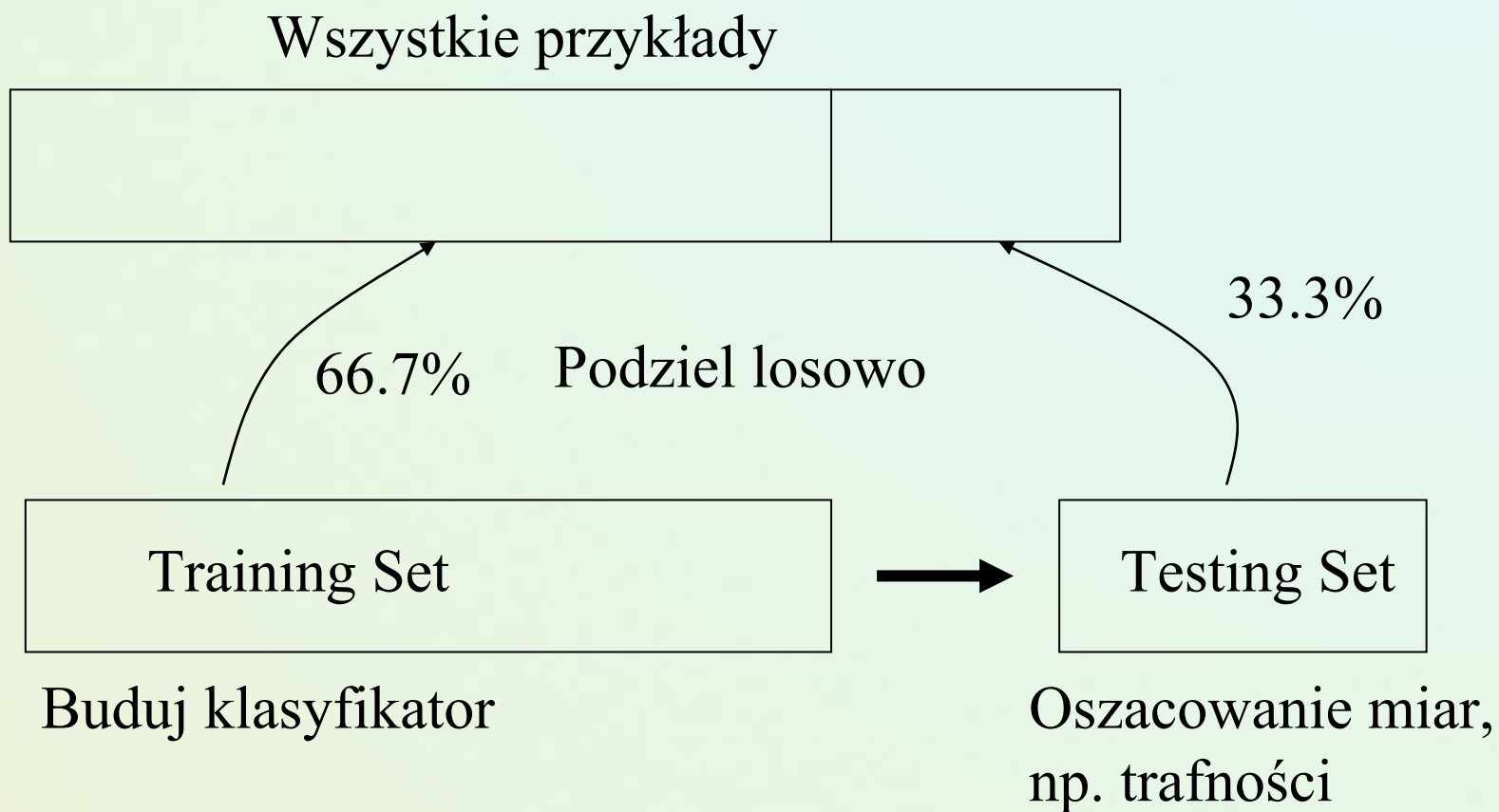


Empiryczne metody estymacji

- **Techniki podziału: „hold-out”**
 - Użyj dwóch niezależnych zbiorów: uczącego (2/3), testowego (1/3)
 - Jednokrotny podział losowy stosuje się dla dużych zbiorów (hold-out)
- **„Cross-validation” - Ocena krzyżowa**
 - Podziel losowo dane w k podzbiorów (równomierne lub warstwowe)
 - Użyj $k-1$ podzbiorów jako części uczącej i pozostałej jako testującej (k -fold cross-validation).
 - Oblicz wynik średni.
 - Stosowane dla danych o średnich rozmiarach (najczęściej $k = 10$)
Uwaga opcja losowania warstwowego (ang. stratified sampling).
- **Leaving-one-out**
 - Dla małych rozmiarów danych.
 - „Leaving-one-out” jest szczególnym przypadkiem, dla którego liczba iteracji jest równa liczbie przykładów

Jednokrotny podział (hold-out)

– duża liczba przykładów (> tysiące)



Inne metody klasyfikacyjne

- sztuczne sieci neuronowe
- klasyfikacja bayesowska
- analiza dyskryminacyjna (statystyczna)
- metody k-najbliższych sąsiadów
- metoda wektorów wspierających (SVM)
- algorytmy genetyczne
- metody oparte na logice matematycznej (ILP)
- ...

Klasyfikacja Bayesowska

- Probabilistic learning: Pozwala na obliczanie prawdopodobieństw związanych z hipotezami; efektywne podejście do wielu problemów praktycznych.
- Incremental: Każdy nowy przykład może zmienić oszacowanie prawdopodobieństw; Możliwość uwzględniania prawdopodobieństw a’piori.
- Predykcja probabilistyczna: Predykcja wielu hipotez, “ważonych” za pomocą prawdopodobieństw.
- Standard: Pomimo, że w pewnych przypadkach metody bayesowskie są nieatrakcyjne obliczeniowo, to mogą dostarczać standard optymalnych decyzji, z którym można porównywać inne metody.

Twierdzenie Bayes'a

- Dla zbioru uczącego D , *prawdopodobieństwo posteriori hipotezy h* , $P(h|D)$ wynika z twierdzenia Bayes'a:

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

- MAP (maximum posteriori) hypothesis:

$$h_{MAP} \equiv \operatorname{argmax}_{h \in H} P(h|D) = \operatorname{argmax}_{h \in H} P(D|h)P(h).$$

- Praktyka: założenie znajomościprawdo-podobieństw *a priori*, duże koszty obliczeniowe.
- Przykład: System d'Dombal → spójrz na dodatkowe slajdy.

Zbiór uczący

- Przykład
Quinlan'a
(*Play a
game*).

Outlook	Tempreature	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

Naiwny klasyfikator bayesowski (II)

- Faza uczenia → obliczanie prawdopodobieństw warunkowych

Outlook	P	N		Humidity	P	N
sunny	2/9	3/5		high	3/9	4/5
overcast	4/9	0		normal	6/9	1/5
rain	3/9	2/5				
Temperature				Windy		
hot	2/9	2/5		true	3/9	3/5
mild	4/9	2/5		false	6/9	2/5
cool	3/9	1/5				

Analiza nowego przypadku

Outlook			Temperature			Humidity			Windy			Play	
Yes	No		Yes	No		Yes	No		Yes	No	Yes	No	
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/14	5/14
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5		
Rainy	3/9	2/5	Cool	3/9	1/5								

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

- Nowa sytuacja

Likelihood of the two classes

$$\text{For "yes"} = 2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0053$$

$$\text{For "no"} = 3/5 \times 1/5 \times 4/5 \times 3/5 \times 5/14 = 0.0206$$

Conversion into a probability by normalization:

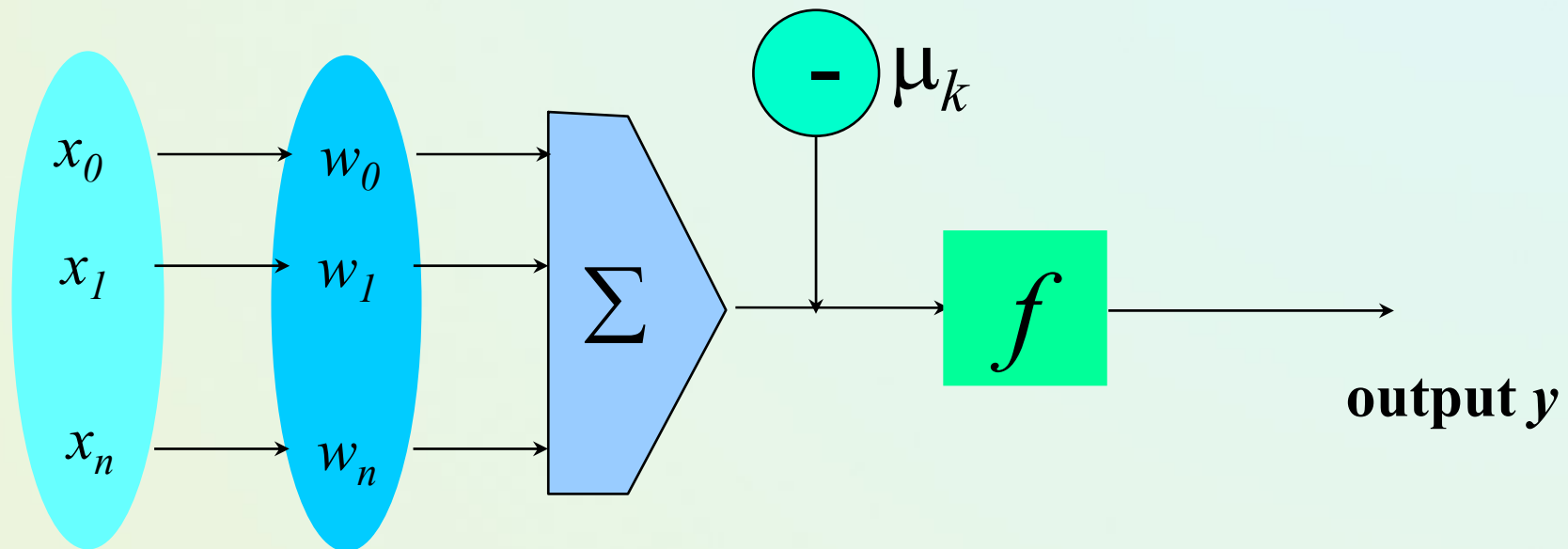
$$P(\text{"yes"}) = 0.0053 / (0.0053 + 0.0206) = 0.205$$

$$P(\text{"no"}) = 0.0206 / (0.0053 + 0.0206) = 0.795$$

Sieci Neuronowe

- Zalety
 - Często wysoka trafność klasyfikacji
 - Odporne na „zaszumione” dane
 - Wyjście może być dyskretne, liczbą rzeczywistą, lub wektorem dyskretnym / liczbowym
 - Po nauczaniu, względnie szybki klasyfikator
- Krytycyzm
 - Długi i złożony proces uczenia
 - Pozyskana wiedza niemożliwa do wyjaśnienia (wagi).
 - Trudno uwzględniać wiedzę dziedzinową

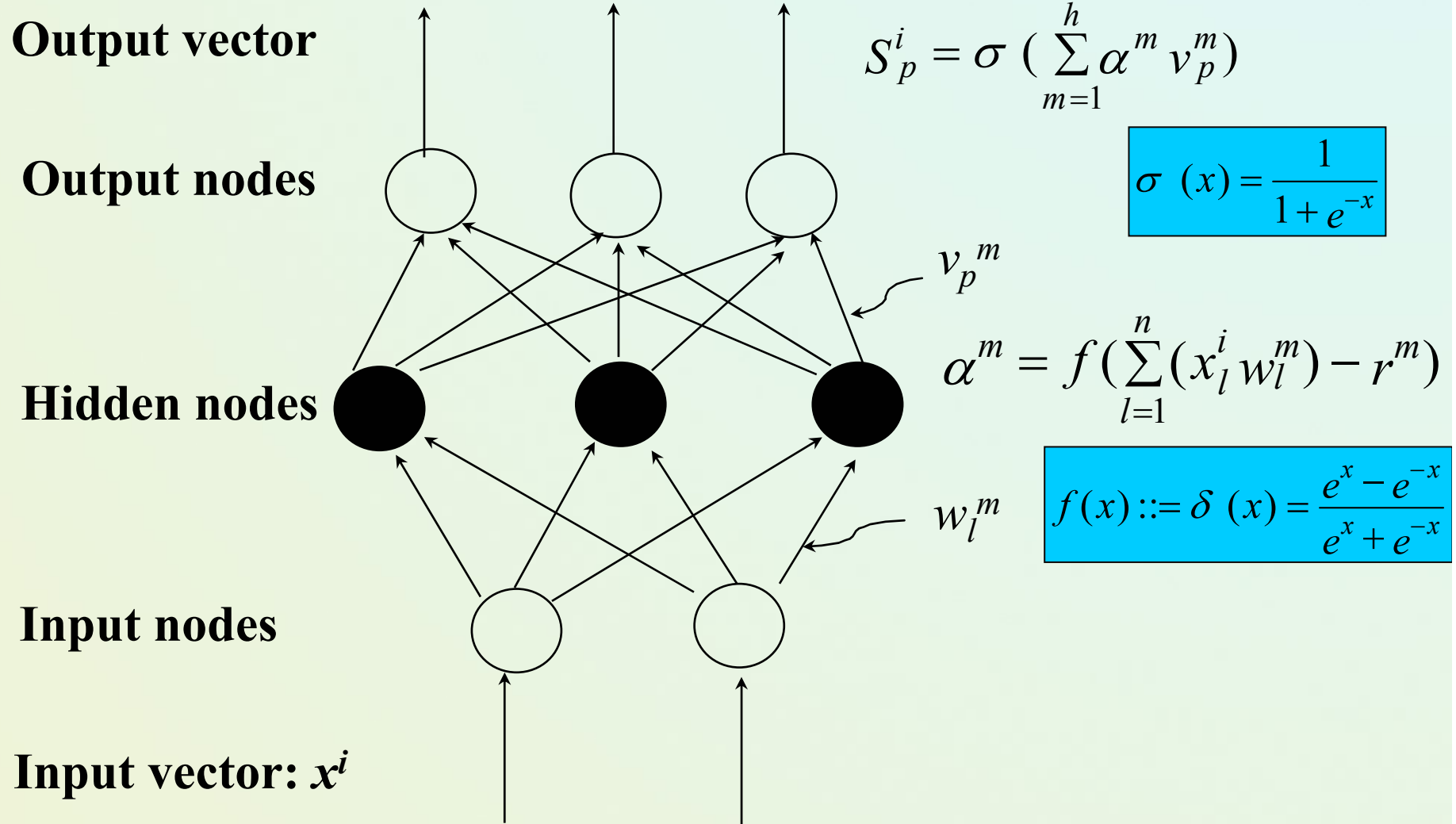
Sztuczny neuron



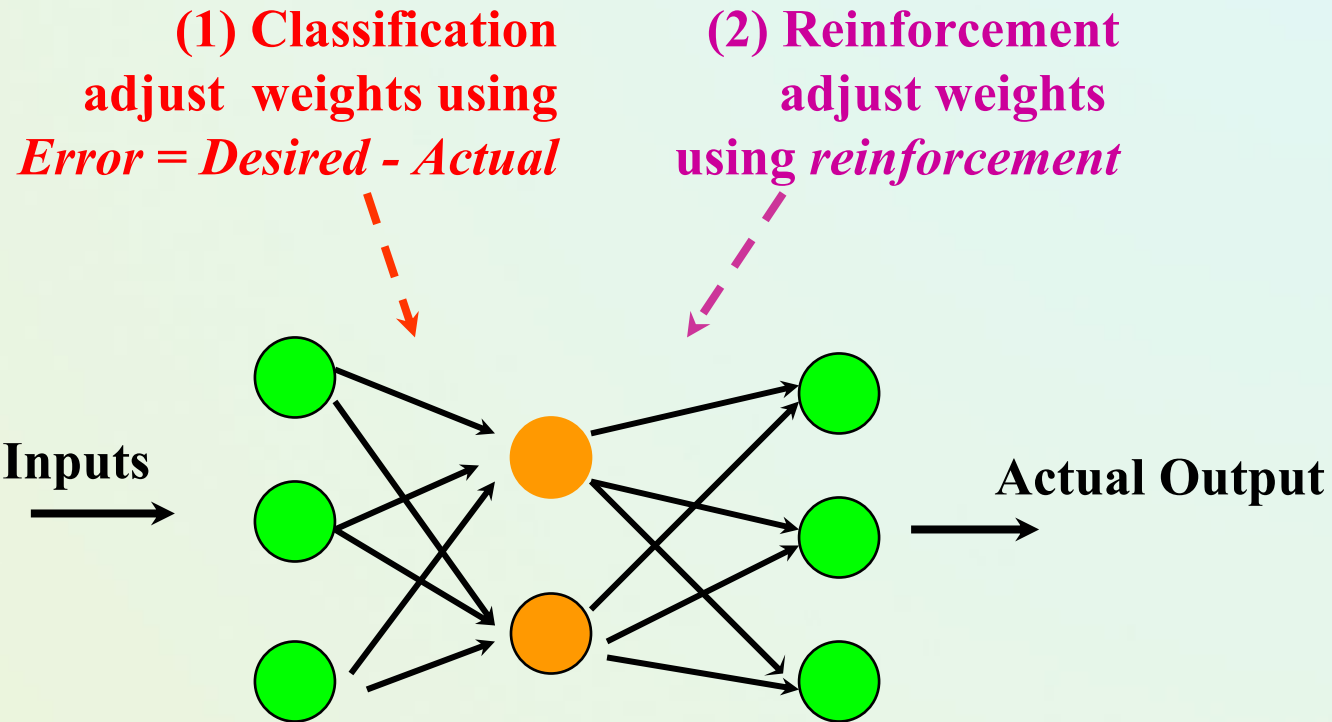
Input **weight** **weighted** **Activation**
vector x **vector w** **sum** **function**

- n -wymiarowy wektor x jest odwzorowywany w zmienną y za pomocą iloczynu skalarnego i nieliniowej funkcji przejścia

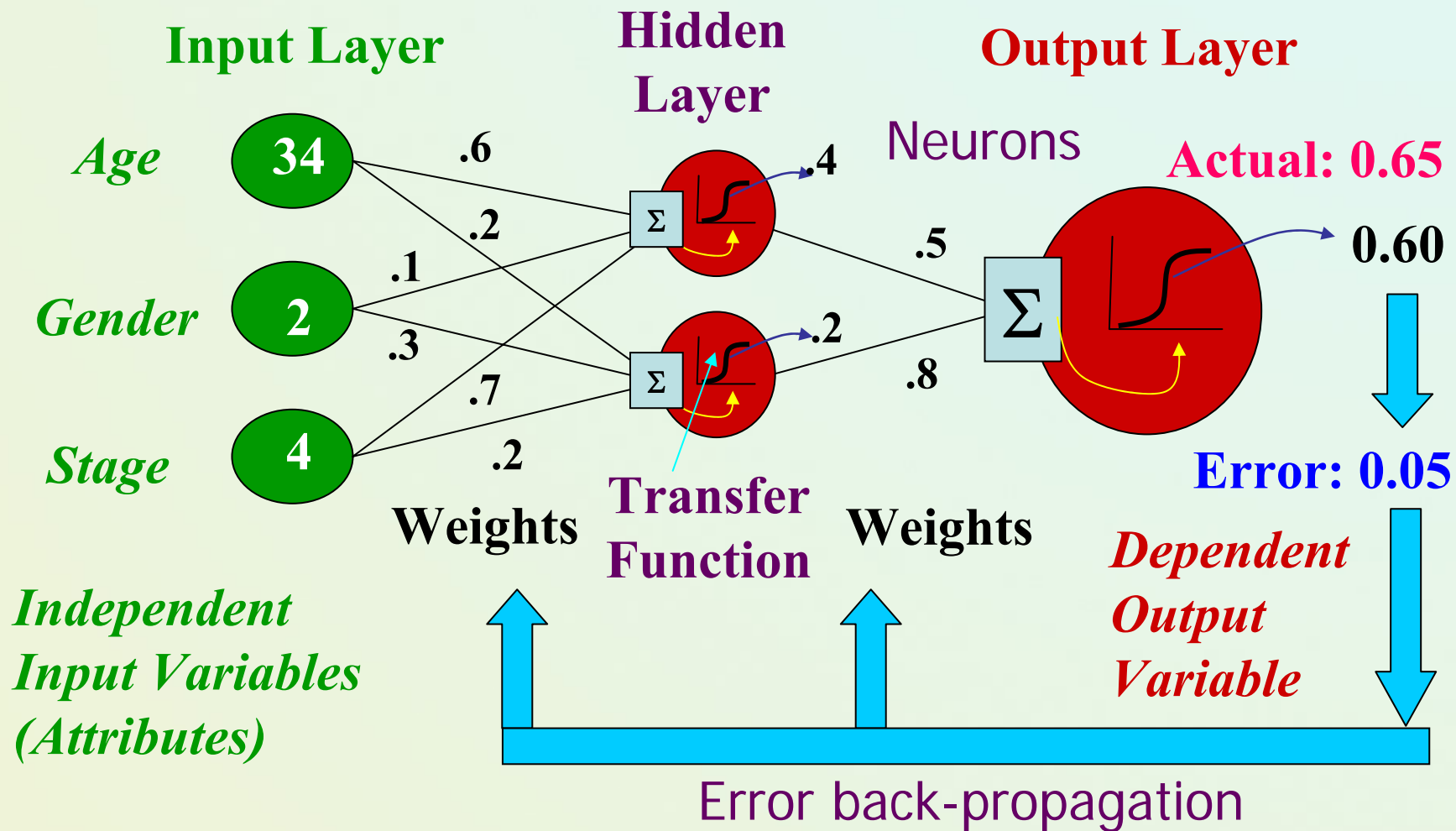
Multi Layer Perceptron



Paradygmaty Uczenia się (I) - Sieć Warstwowa



The Neural Network (NN) Approach



STATISTICA

Create Network [?] [X]

Type: **Multilayer Perceptron** [Advise]

Time Series [Create]

Steps: **1** [Lookahead: **0**] [Close]

Pre/Post Processing

Inputs: **4** [Outputs: **1**] [No Layers: **3**]

Convert	Units
LENGTH	Minimax
WIDTH	Minimax
PLENGTH	Minimax
PWIDTH	Minimax
FLOWER	One-of-N

Units	
Layer 1	4
Layer 2	3
Layer 3	3

Back Propagation [?] [X]

Epochs: **100** [Train]

Learning rate: **0.6** [Reinitialize]

Momentum: **0.3** [Log Weights]

Noise: **0** [Stop]

Shuffle Cases

Cross verification

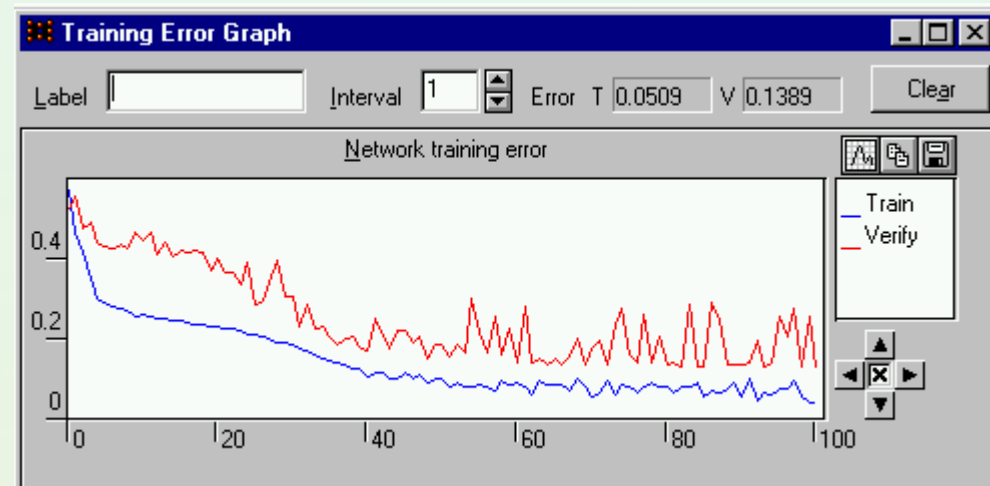
[Close]

Run Data Set [?] [X]

Outputs shown: **Variables** [Run]

RMS Error Train: **0.06435** Verify: **0** Test: **0**

RESULT	T. RESULT	E. RESULT	Error	
01	0.06583	0	0.06583	0.06583
fred	0.9384	1	-0.0616	0.0616
03	0.9376	1	-0.0624	0.0624
04	0.0674	0	0.0674	0.0674



- Tworzenie i uruchomienie sieci

Przykład pakietu SNNS

The screenshot displays the Stuttgart Neural Network Simulator (SNNS) interface, which is divided into several windows:

- snn-manager:** A menu-driven control window with options like FILE, CONTROL, INFD, DISPLAY, 3D DISPLAY, GRAPH, BIGNET, PRUNING, CRSCODE, KOHONEN, WEIGHTS, PROJECTION, ANALYZER, DIMENSION, PRINT, HELP, and RPC. It also shows coordinates (x:15, y:0) and a status bar.
- banner:** Displays the title "SNNS V4.0" and "PART NEURAL NETWORK SIMULATOR".
- 3D-display:** A 3D visualization of a neural network with nodes and connections.
- RPCSETUP:** A configuration window for connection types (UDP, TCP) and kernel types (KERNEL, COOP-MASTER). It includes settings for timeouts and switching between local and remote execution.
- RPCSELECTITEMLIST:** A window for selecting information to show in the RPC panel, with fields for INFO and SELECTED.
- SNNS RPC PANEL:** A window showing the current host and a table of active processes.
- Terminal (matisse):** A terminal window showing the command `dm -program 572352910 -version 40 -uid 20139 < /dev/null > /dev/null 2>&1 &''&<-` and the output `New Kernel on matisse ,PID 2122 added`.

No.	HOSTNAME	KERNELNO	STATUS	SSE
0	direct local	0	idle	0,52159
1	vazarnly	400020139	idle	0,00000
2	wondrian	400020139	idle	0,00000
3	matisse	400020139	idle	0,00000

- *Stuttgart Neural Network Simulator*

Przykład: Klasyfikacja elektroencefalogramu

- **Zadanie:** klasyfikacja elektroencefalogramu do jednej z trzech klas:
 1. normal (*norm*)
 2. schizophrenic (*shiz*)
 3. obsessive compulsive disorder (*ocd*)

Dane wejściowe:

- postać oryginalna danych: wynik badania elektroencefalograficznego zaklasyfikowany przez neurologów do jednej z trzech klas: 19 kanałów próbkowanych z częstotliwości 128 Hz, pacjent w spoczynku

Przetwarzanie wstępne:

- usunięcie artefaktów (np. zakłóceń wynikłych z poruszeń pacjenta).
- usunięcie składowych o wysokich częstotliwościach (np. szumów) poprzez obróbkę filtrem dolnoprzepustowym (30 Hz),...
- z 19 sygnałów do dalszej obróbki wybrany kanał Cz (pobierany ze szczytu głowy) → wcześniejsze badania wykazały największą przydatność do celów diagnostycznych (w przybliżeniu sygnał dostarczany tym kanałem jest liniową kombinacją części innych sygnałów) .
- Wybór okna czasowego (pierwsze 250 sekund) – stan największego uspokojenia pacjenta
-

Klasyfikacja elektroencefalogramu

- **Ekstrakcja cech:** Metoda: autoregresji (AR); polega na przybliżeniu przebiegu przez rekurencyjną zależność:

$$s_t = \sum_{j=1}^N \alpha_j s_{t-j} + \varepsilon_t$$

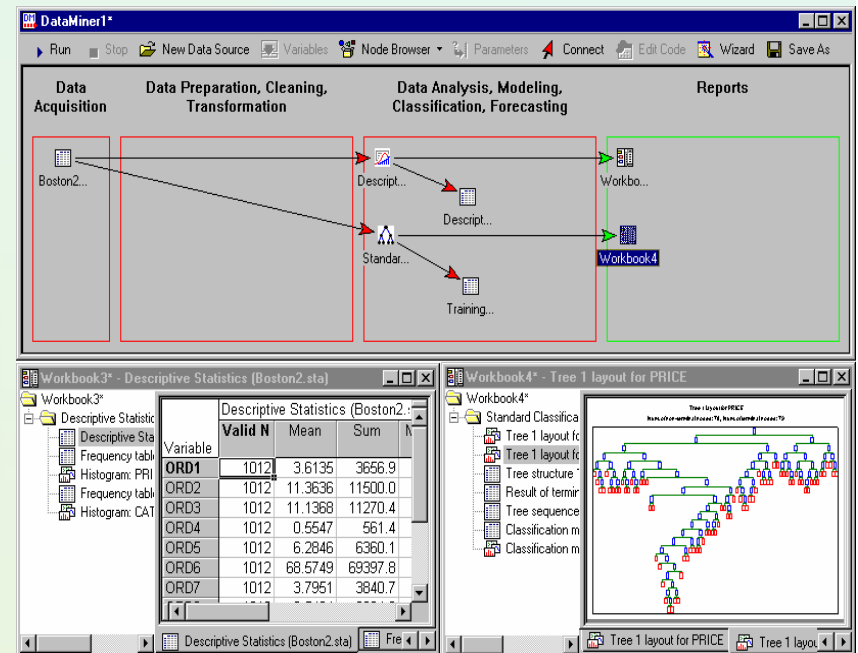
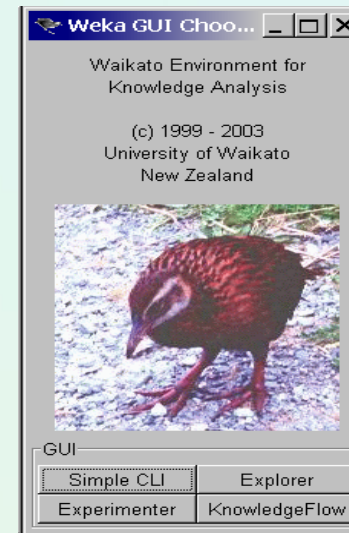
- s_t - sygnał w chwili t , ε pewna zmienna losowa o rozkładzie normalnym
- Szukanymi parametrami opisującymi przebieg są α_j dla $j=1..N$; oblicza się specjalnymi metodami - są atrybutami podawanymi na wejścia sieci. Współczynniki mają podobną interpretację jak elementy widma częstotliwości otrzymanego np. przy pomocy FFT, lecz są bardziej odporne na szумы.

EEG – konstrukcja sieci neuronowej

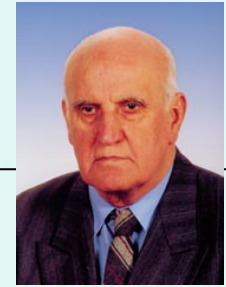
- Topologia: dwuwarstwowy perceptron o architekturze 8-15-3 (8 neuronów w warstwie wejściowej, 15 w warstwie ukrytej i 3 w wyjściowej)
- Algorytm uczący: klasyczny *backpropagation*
- Trafność klasyfikowania: 91%
 - Więcej w: Tsoi A.C., So D.S.C, Sergejew A.: Classification of Electroencephalogram using Artificial Neural Networks. Technical Report, University of Queensland, Australia, March 1993
 - Przejrzyj inne książki, aby poznać więcej zastosowań, np. R.Tadeusiewicza.

Dostępne narzędzia, implementacje metod

- Open source lub freeware
 - **WEKA** – Waikato
 - MOA – Datastreams Waikato
 - Rapid Miner (YALE)
 - Orange Lubljana
 - KNIME Konstanz
 - MLC++ Stanford
 - SSNN Stutgart
 - **ROSE** i inne (PP Poznań)
- Systemy komercyjne
 - SAS Institute: Enterprise Miner
 - Statistica Data Miner
 - SPSS : Clementine
 - Oracle 9i Miner
 - IBM: QUEST and Intelligent Miner
 - Silicon Graphics: MineSet
 - Review 10 Top Data Mining Tools Edler

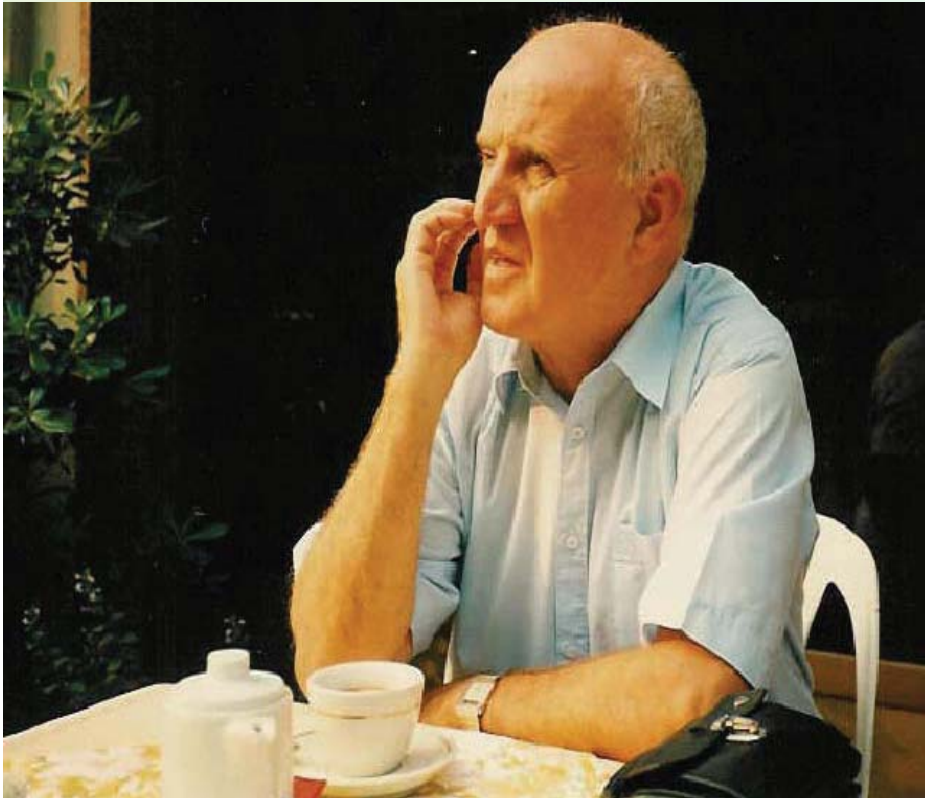


Rough Sets – Zbiory przybliżone



- Zainteresowanie niedoskonałościami i niespójnościami w danych.
- - **Inconsistency – niespójność w opisie obiektów** nie zawsze jest wynikiem błędów lub „szumu” informacyjnego
 - W niektórych zastosowaniach pożądana jest identyfikacja pewnej części wiedzy pozyskanej z dostępnych danych i oddzielenie jej od części możliwej lub niepewnej.
 - Istnieją różne typy i semantyki niespójności w danych.
- **Rough sets theory** – teoria zbiorów przybliżonych (Z.Pawlak) i jej uogólnienia (np. R.Słowiński, A.Skowron, Y.Yao, S.Greco,...) dostarcza podstaw teoretycznych do analizy niespójnych danych.

Zdzisław I. Pawlak (1926-2006)



- Najbardziej znany na świecie polski naukowiec z zakresu informatyki
- Brał udział w konstrukcji polskich komputerów (latach 50-60)
- Pierwsze prace polskie prace naukowe publikowane w USA (1953)
- Logiczne podstawy informatyki
- Teoria automatów
- Maszyny bezadresowe
- Języki wyszukiwania informacji
- Teoria zbiorów przybliżonych

Rough Sets (2)

- Opiera się na badaniu relacji zachodzących między opisami obiektów.
- Podstawowe relacje: nierozróżnialność, podobieństwo, dominacja
- Jeśli opisy obiektów za pomocą atrybutów są niespójne (sprzeczne), to tworzy się tzw. przybliżenia zbiorów (klas obiektów).
- **Rough set** – zbiór przybliżonych jest **parą przybliżeń** zbioru!
- Może być użyta dla oceny przydatności atrybutów dla przybliżenia klasyfikacji, selekcji atrybutów, redukcji tablicy danych, zależności między zbiorami atrybutów, generacji wzorców z danych (reguł, templates, reguł asocjacyjnych), itd.
- Oraz wiele innych ...

Przybliżenie klasyfikacji obiektów

<i>U</i>	<i>Headache</i>	<i>Temp.</i>	<i>Flu</i>
<i>U1</i>	Yes	Normal	No
<i>U2</i>	Yes	High	Yes
<i>U3</i>	Yes	Very-high	Yes
<i>U4</i>	No	Normal	No
<i>U5</i>	<i>No</i>	<i>High</i>	<i>No</i>
<i>U6</i>	<i>No</i>	<i>Very-high</i>	<i>Yes</i>
<i>U7</i>	<i>No</i>	<i>High</i>	<i>Yes</i>
<i>U8</i>	<i>No</i>	<i>Very-high</i>	<i>No</i>

Klasy relacji nierozróżnialności (zbiory elementarne) dla zbioru atrybutów $R = \{Headache, Temp.\}$:
 $\{u1\}, \{u2\}, \{u3\}, \{u4\}, \{u5, u7\}, \{u6, u8\}$.

$$X1 = \{u \mid Flu(u) = yes\}$$
$$= \{u2, u3, u6, u7\}$$

$$\underline{RX1} = \{u2, u3\}$$

$$\overline{RX1} = \{u2, u3, u6, u7, u8, u5\}$$

$$X2 = \{u \mid Flu(u) = no\}$$

$$= \{u1, u4, u5, u8\}$$

$$\underline{RX2} = \{u1, u4\}$$

$$\overline{RX2} = \{u1, u4, u5, u8, u7, u6\}$$

Przybliżenie klas decyzyjnych $X1$ oraz $X2$

$$R = \{Headache, Temp.\}$$

$$U/R = \{ \{u1\}, \{u2\}, \{u3\}, \{u4\}, \{u5, u7\}, \{u6, u8\} \}$$

$$X1 = \{u \mid Flu(u) = yes\} = \{u2, u3, u6, u7\}$$

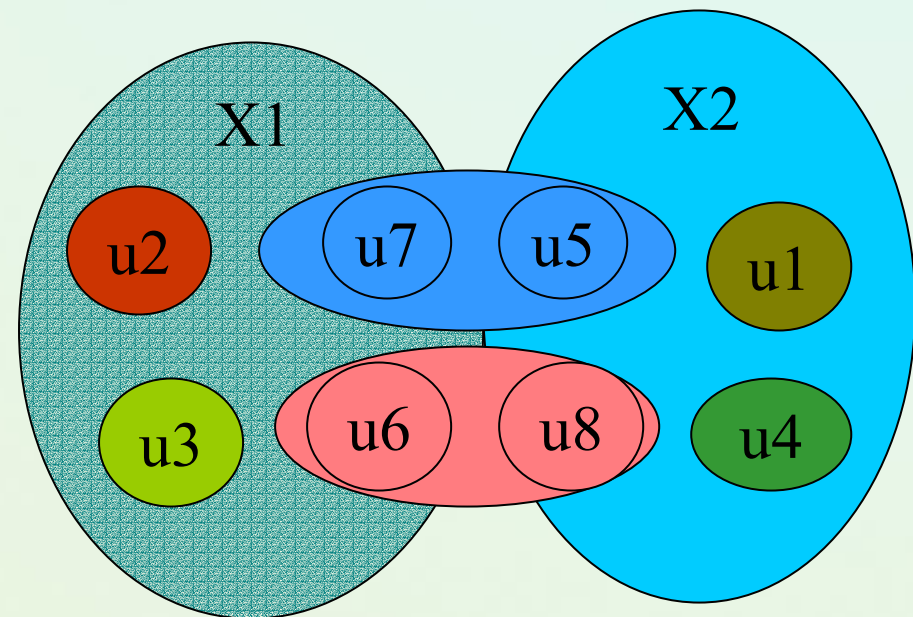
$$X2 = \{u \mid Flu(u) = no\} = \{u1, u4, u5, u8\}$$

$$\underline{RX1} = \{u2, u3\}$$

$$\overline{RX1} = \{u2, u3, u6, u7, u8, u5\}$$

$$\underline{RX2} = \{u1, u4\}$$

$$\overline{RX2} = \{u1, u4, u5, u8, u7, u6\}$$



Przykład redukcji tablicy decyzyjnej

<i>U</i>	<i>Headache</i>	<i>Muscle pain</i>	<i>Temp.</i>	<i>Flu</i>
<i>U1</i>	Yes	Yes	Normal	No
<i>U2</i>	Yes	Yes	High	Yes
<i>U3</i>	Yes	Yes	Very-high	Yes
<i>U4</i>	No	Yes	Normal	No
<i>U5</i>	No	No	High	No
<i>U6</i>	No	Yes	Very-high	Yes

$$\begin{aligned}
 CORE &= \{Headache, Temp\} \cap \\
 &\quad \{MusclePain, Temp\} \\
 &= \{Temp\}
 \end{aligned}$$

Reduct1 = {Muscle-pain, Temp.}



<i>U</i>	<i>Muscle pain</i>	<i>Temp.</i>	<i>Flu</i>
<i>U1, U4</i>	Yes	Normal	No
<i>U2</i>	Yes	High	Yes
<i>U3, U6</i>	Yes	Very-high	Yes
<i>U5</i>	No	High	No

Reduct2 = {Headache, Temp.}



<i>U</i>	<i>Headache</i>	<i>Temp.</i>	<i>Flu</i>
<i>U1</i>	Yes	Normal	No
<i>U2</i>	Yes	High	Yes
<i>U3</i>	Yes	Very-high	Yes
<i>U4</i>	No	Normal	No
<i>U5</i>	No	High	No
<i>U6</i>	No	Very-high	Yes

LEM2 – An Example (1)

<i>U</i>	<i>Headache</i>	<i>Nausea</i>	<i>Temp.</i>	<i>Flu</i>
<i>x1</i>	no	no	normal	No
<i>x2</i>	yes	no	high	Yes
<i>x3</i>	yes	yes	high	Yes
<i>x4</i>	yes	no	normal	No
<i>x5</i>	no	no	high	No
<i>x6</i>	no	no	high	Yes

IND: {*x1*}, {*x2*}, {*x3*}, {*x4*}, {*x5,x6*}

YES: lower appr. {*x2,x3*}

upper {*x2,x3,x5,x6*}

NO: lower approx. {*x1,x4*}

upper {*x1,x4,x5,x6*}

Reguły (pewne) dla dolnego przybliżenia klasy (**Flue=Yes**) – przykłady
pozytywne {*x2,x3*}

(headache,yes)	{ <i>x2,x3</i> ⁺ ; <i>x4</i> ⁻ }
(nausea,no)	{ <i>x2</i> ⁺ ; <i>x1,x4,x5,x6</i> ⁻ }
(nausea,yes)	{ <i>x3</i> ⁺ }
(temperature,high)	{ <i>x2,x3</i> ⁺ ; <i>x5,x6</i> ⁻ }

Wybierz *t1* (headache,yes), ale {*x2,x3*⁺ ; *x4*⁻} $\not\subseteq$ {*x2,x3*}, poszukuj
następnego warunku; wybierz (temperature,high),

Teraz $t1 \cap t2 = \{x2,x3^+ ; x4^-\} \cap \{x2,x3^+ ; x5,x6^-\} \subseteq \{x2,x3\}$

reguła (**headache=yes**) \cap (**temperature=high**) \rightarrow (**Flue=Yes**) pokrywa
wszystkie przykłady pozytywne.

LEM2 – An Example (2)

<i>U</i>	<i>Headache</i>	<i>Nausea</i>	<i>Temp.</i>	<i>Flu</i>
<i>x1</i>	no	no	normal	No
<i>x2</i>	yes	no	high	Yes
<i>x3</i>	yes	yes	high	Yes
<i>x4</i>	yes	no	normal	No
<i>x5</i>	no	no	high	No
<i>x6</i>	no	no	high	Yes

IND: {x1}, {x2}, {x3}, {x4}, {x5,x6}

YES: lower appr. {x2,x3}

upper {x2,x3,x5,x6}

NO: lower approx. {x1,x4}

upper {x1,x4,x5,x6}

Reguły pewne (**Flue=No**): Dolne przybliżenie {x1,x4}

(headache,no) {x1+; x5,x6-}

(headache,yes) {x4+ ; x2,x3-}

(nausea,no) {x1,x4+;x2,x5,x6-}

(temperature,normal) {x1,x4+ ; ∅}

Wybierz t_1 (temperature,normal),

$t_1 = \{x1,x4+ ; \emptyset-\} \subseteq \{x1,x4\}$

Reguła (**temperature=normal**) \rightarrow (**Flue=No**) pokrywa wszystkie przykłady pozytywne.

Przykładowe zastosowania teorii zbiorów przybliżonych w analizie medycznych danych.

- Klasyfikacja pacjentów z chorobą wrzodową dwunastnicy leczonych metodą wysoce wybiórczej wagoatomii.
- Wieloetapowe leczenie ostrego zapalenia trzustki płukaniem otrzewnej.
- Leczenie kamicy moczowej za pomocą ESWL.
- Analiza leczenia pacjentów z obrażeniami wielo-narządowymi przyjmowanych w izbie przyjęć.
- Weryfikacja diagnozowania uszkodzeń stawu kolanowego na podstawie MRI.
- Rokowanie przeżycia przy leczeniu raka piersi.
- Wspomaganie decyzji w przypadku bólów brzucha u dzieci w izbie przyjęć.
- ...

- Więcej w: Z.Pawlak, K.Słowiński, J.Stefanowski: Teoria zbiorów przybliżonych w analizie danych medycznych, 2002.



Leczenie pacjentów z chorobą wrzodową dwunastnicy

- **Celem** było określenie wskazań do przeprowadzania zabiegu chirurgicznego wysoce wybiórczej wagoatomii (*HSV*) dla pacjentów cierpiących na chorobę wrzodową dwunastnicy.
- **Dane**: Pacjenci opisani są za pomocą 11 atrybutów i przydzieleni do 4 klas wyrażających skuteczność zabiegu chirurgicznego. Klasy w sensie liczności są niezrównoważone (silna przewaga grupy wyleczonych pacjentów). Atrybuty są różnego typu (nominalnego, porządkowego jak i liczbowego).

Leczenie pacjentów z chorobą wrzodową dwunastnicy

Aspekty metodologiczne:

- Postawione zadania:
 - ocena znaczenia atrybutów i ich podzbiorów dla przybliżenia klasyfikacji pacjentów;
 - poszukiwanie tzw. modeli pacjentów charakterystycznych dla klas decyzyjnych;
 - poszukiwanie reprezentacji zależności pomiędzy wartościami wybranych atrybutów warunkowych a decyzyjnym w postaci reguł decyzyjnych.
- Zastosowano teorię zbiorów przybliżonych oraz metody indukcji reguł.
- Na podstawie rezultatów analizy danych oraz wiedzy o leczeniu sformułowano algorytmy kliniczne.

Analiza MRI zdjęć uszkodzeń stawu kolanowego

- **Sformułowanie problemu:** Problem dotyczył weryfikacji klinicznej diagnoz pewnej formy uszkodzenia stawu kolanowego (*anterior cruciate ligament* – ACL) postawionych na podstawie symptomów zaobserwowanych na zdjęciach wykonanych techniką magnetycznego rezonansu jądrowego (MRI). Procedura chirurgiczna artroskopii umożliwia bezpośrednią weryfikację uszkodzenia stawu, lecz jest techniką inwazyjną i wykonana niepoprawnie może prowadzić do komplikacji.

Dane i ich charakterystyka:

- Wybrano grupę 140 pacjentów, dla których wykonano zarówno standardowe procedury diagnostyczne i badanie za pomocą MRI jak i artroskopię. Dla 100 pacjentów wykluczono, a dla 40 potwierdzono występowanie uszkodzenia ACL.
- Przy współpracy z ekspertami medycznymi oraz na podstawie wstępnej analizy statystycznej wyselekcjonowano zbiór 6 atrybutów opisujących pacjentów. Większość atrybutów wynikała z pomiarów przeprowadzonych na zdjęciu MRI (PCL index; pomiary X i Y). Pozostałe (wiek, płeć, strona ciała)

Przetwarzanie wstępne - dyskretyzacja

Table 1. Comparison of various discretizations

Evaluation measures	<i>Expert</i>	<i>MI</i>	<i>ME</i>	<i>MD</i>
Rough sets results:				
accuracy of class 1	0.653	0.822	0.951	1.0
accuracy of class 2	0.843	0.922	0.980	1.0
quality of classification	0.879	0.943	0.987	1.0
no. of attributes in a core	6	6	4	4
number of reducts	1	1	2	1
LEM2 results:				
number of rules	23	20	16	20
total classification accuracy	85%	88.86%	92.86%	90%
sensitivity of Class 1	67.2%	80%	90%	87.5%
average rule strength (obj.)	9.65	8.57	10.62	4.93
average rule length	3.67	2.67	2.47	2.73

AGE (years): [min,16.5), [16.5,35), [35,max]

X (mm): [min,8.5), [8.5,11.75), [11.75,14.5), [14.5,max]

Y (mm): [min,2.75), [2.75,3.75), [3.75,4.75) [4.75,max]

PCLINDEX: [min,3.225), [3.225,3.71), [3.71,4.125), [4.125,4.535), [4.535,max]

Przybliżenia klasyfikacji

- Jakość przybliżenia = 0.9429

Table 2. Rough approximations of the patients' classification

class	number of patients	cardinality of lower approx.	cardinality of upper approx.	accuracy of class
1	40	39	41	0.9512
2	100	99	101	0.9802

Znaczenie atrybutów

Table 3. Significance of attributes of approximating the patients' classification

Attribute	AGE	SEX	SIDE	X	Y	PCLINDEX
Significance	0.057	0.043	0.021	0.00	0.00	0.071

- Rdzeń:
- PCLINDEX, AGE, SEX, SIDE

Minimalny zbiór reguł

- MODLEM

rule 1. *if* (PCLINDEX < 3.225) *then* Class1 [26, 65%]

rule 2. *if* (AGE=[16.5,35])^(PCLINDEX=[3.225,3.71)) *then* Class1 [6, 15%]

rule 3. *if* (SEX=MALE) ^ (SIDE=RIGHT) ^ (PCLINDEX=[3.225,3.71)) *then* Class1 [3, 7.5%]

rule 4. *if* (AGE=[16.5,35]) ^ (PCLINDEX=[3.71,4.125)) ^ (X≥14.5) *then* Class1 [2, 5%]

rule 5. *if* (X=[8.5,11.75)) ^ (PCLINDEX=[4.125,4.535)) ^ (SEX=MALE) *then* Class1 [1, 2.5%]

rule 6. *if* (X=[8.5,11.75)) ^ (PCLINDEX=[3.225,3.71)) ^ (AGE≥ 35) *then* Class1 [2, 5%]

rule 7. *if* (PCLINDEX=[3.71,4.125)) ^ (X=[8.5,11.75)) ^ (SEX=1) *then* Class1 [1, 2.5%]

rule 8. *if* (PCLINDEX≥4.535) *then* Class2 [75, 75%]

rule 9. *if* (SEX=FEMALE) ^ (PCLINDEX=[4.125,4.535)) *then* Class2 [10,10%]

rule 10. *if* (PCLINDEX=[3.71,4.125)) ^ (AGE≥ 35) *then* Class2 [6,6%]

rule 11. *if* (X=[11.75,14.5)) ^ (Y=[2.75,3.75)) ^ (SEX=FEMALE) *then* Class2 [8, 8%]

rule 12. *if* (SIDE=LEFT) ^ (X=[11.75,14.5)) ^ (Y=[2.75,3.75)) *then* Class2 [7, 7%]

rule 13. *if* (PCLINDEX=[3.225,3.71)) ^ (AGE≥ 35) ^ (SEX=MALE) *then* Class2 [2, 2%]

rule 14. *if* (AGE<16.5) *then* Class2 [14, 14%]

rule 15. *if* (PCLINDEX=[3.225,3.71))^(Y=[3.75,4.75))^(AGE≥ 35)^(SIDE=LEFT) *then* Class2 [1,1%]

Satysfakcjonujący zbiór reguł (EXPLORE)

- Rule. Strength $\geq 10\%$ Klasy

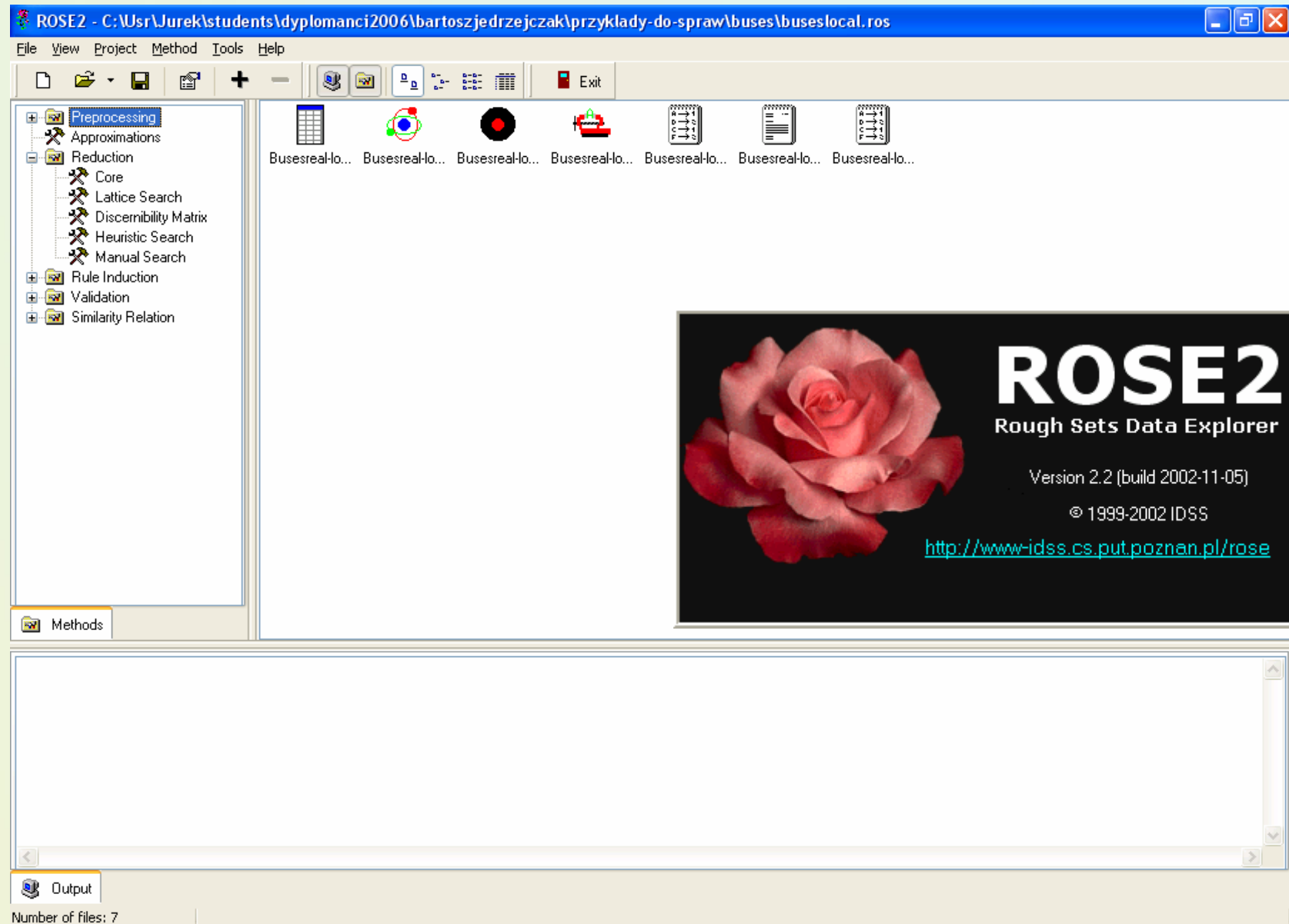
rule 1. *if*(X < 8.5) *then* Class1 [4, 10%]
rule 2. *if*(PCLINDEX < 3.225) *then* Class1 [26, 65%]
rule 3. *if*(AGE=[16.5,35]) \wedge (Y \geq 4.75) *then* Class1 [6, 15%]
rule 4. *if*(AGE=[16.5,35]) \wedge (PCLINDEX = [3.225,3.71]) *then* Class1 [6, 15%]
rule 5. *if*(SEX = MALE) \wedge (Y \geq 4.75) *then* Class1 [9, 22.5%]
rule 6. *if*(SIDE = RIGHT) \wedge (Y \geq 4.75) *then* Class1 [7, 17.5%]
rule 7. *if*(X = [8.5,11.75]) \wedge (Y = [3.75,4.75]) *then* Class1 [9, 22.5%]
rule 8. *if*(X = [11.75,14.5]) \wedge (Y \geq 4.75) *then* Class1 [5, 12.5%]
rule 9. *if*(AGE=[16.5,35]) \wedge (X = [8.5,11.75]) \wedge (Y = [2.75,3.75]) *then* Class1 [5, 12.5%]
rule 10. *if*(AGE=[16.5,35]) \wedge (X = [11.75,14.5]) \wedge (Y = [3.75,4.75]) *then* Class1 [4, 10%]
rule 11. *if*(SEX = MALE) \wedge (X = [8.5,11.75]) \wedge (Y = [2.75,3.75]) *then* Class1 [6, 15%]
rule 12. *if*(SEX = MALE) \wedge (X = [8.5,11.75]) \wedge (PCLINDEX = [3.225,3.71]) *then* Class1 [4, 10%]
rule 13. *if*(SEX = MALE) \wedge (Y = [2.75,3.75]) \wedge (PCLINDEX = [3.225,3.71]) *then* Class1 [4, 10%]
rule 14. *if*(SIDE = RIGHT) \wedge (X = [8.5,11.75]) \wedge (Y = [2.75,3.75]) *then* Class1 [5, 12.5%]
rule 15. *if*(AGE < 16.5) *then* Class2 [14, 14%]
rule 16. *if*(PCLINDEX \geq 4.535) *then* Class2 [75, 75%]
rule 17. *if*(AGE \geq 35) \wedge (X \geq 14.5) *then* Class2 [16, 16%]
rule 18. *if*(AGE \geq 35) \wedge (Y < 2.75) *then* Class2 [14, 14%]
rule 19. *if*(SEX = FEMALE) \wedge (X \geq 14.5) *then* Class2 [28, 28%]
rule 20. *if*(SEX = FEMALE) \wedge (PCLINDEX = [4.125,4.535]) *then* Class2 [10, 10%]
rule 21. *if*(SIDE = RIGHT) \wedge (Y < 2.75) *then* Class2 [19, 19%]
rule 22. *if*(SIDE = LEFT) \wedge (X \geq 14.5) *then* Class2 [24, 24%]
rule 23. *if*(X \geq 14.5) \wedge (Y = [2.75,3.75]) *then* Class2 [29, 29%]

Ocena zdolności klasyfikacyjnych

Table 4. Classification performance of decision rules

Rule set	Overall	Class 1		Class 2	
	accuracy	sensitivity	specificity	sensitivity	specificity
minimum	92.86%	90%	94%	94%	90%
satisfactory	93.57%	85%	97%	97%	85%

Oprogramowanie do TZP → ROSE (IDSS)



Rose → redukcja

- Redukt vs. significance for classification

The image shows two overlapping windows from a software application. The background window is 'Reduct Viewer' and the foreground window is 'Selecting attributes'.

Reduct Viewer - C:\Usr\Jurek\students\dyplomanci20...

#	Reduct	Length
1	Compr_preasure, blacking, torque	
2	MaxSpeed, oil_cons	
3	MaxSpeed, Compr_preasure	
4	Compr_preasure, oil_cons	
5	Compr_preasure, horsepower	

Number of reducts: 5

Selecting attributes

Chosen attributes:

Attribute name	Quality loss
Compr_preasure	0.382

Removed attributes:

Attribute name	Quality gain
blacking	0.579
horsepower	0.618
MaxSpeed	0.618
oil_cons	0.618
summer_cons	0.039
torque	0.592
winter_cons	0.408

Automatic calculate quality changes

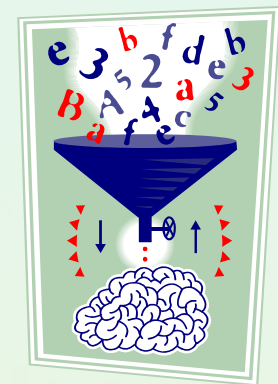
Decision attribute: D1
Classification estimation: Quality

Classification quality for all attributes: 1.000
Classification quality for chosen attributes: 0.382

Buttons: Remove >>, << Add, Back, Add reduct to list, Show list of reducts, Calculate quality changes, Close

Podsumowanie:

- Nowe technologie dostarczają wielu danych
 - data mining pomaga odnaleźć interesujące zależności
- Data Mining i Uczenie Maszynowe
- Wiele prac, ale nie oczekujemy masowych „automatycznych” zastosowań ?
- Podstawowe zadania:
 - Klasyfikacja nadzorowana, grupowanie, ...
- Omówione metody:
 - Drzewa i reguły, klasyfikacja bayesowska, sztuczne sieci neuronowe, teoria zbiorów przybliżonych (mniej ML bardziej ...)



Czy są jakieś pytania?



I to by było na tyle ...

Nie zadawaj się tym
co usłyszałeś – poszukuj więcej!
Czytaj książki oraz samodzielnie
eksploruj dane!

