

# Strojno učenje

Marko Bohanec

Institut Jožef Stefan, Ljubljana in  
Univerza v Novi Gorici

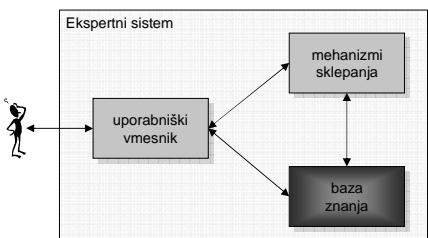
Marko Bohanec

## Kazalo

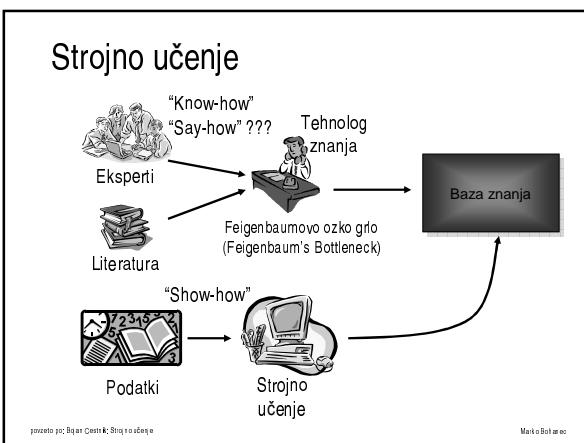
- Uvod
  - Zakaj strojno učenje?
  - Kaj je strojno učenje?
  - Zahteve strojnega učenja
  - Področja uporabe
  - Literatura in viri
- Metode strojnega učenja

Marko Bohanec

## Arhitektura ES



Marko Bohanec




---

---

---

---

---

---

---

---

---

---

**Primer**

Dobiček	Starost	Konkurenca	Vrsta
pada	staro	ne	SW
pada	srednje	da	SW
narašča	srednje	ne	HW
pada	staro	ne	HW
narašča	novo	ne	HW
narašča	novo	ne	SW
narašča	srednje	ne	SW
narašča	novo	da	SW
pada	srednje	da	HW
pada	staro	da	SW

izvor: Bojan Čestrik; Strojno učenje  
 Marko Botičarec

---

---

---

---

---

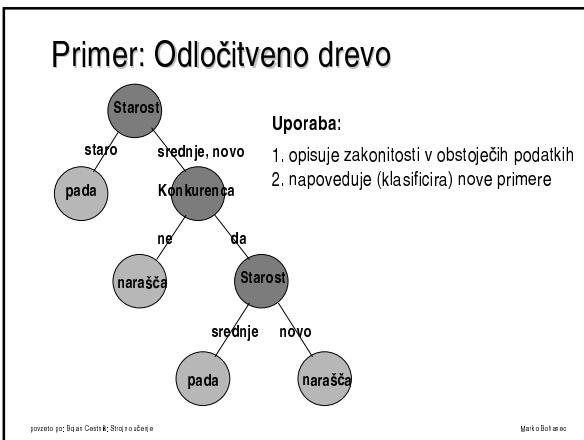
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

## Metode strojnega učenja

- Statistične metode
  - Bayesov klasifikator
  - *k*-najbližjih sosedov (*k*-Nearest Neighbors, *k*-NN)
  - diskriminantna analiza
- Simbolično induktivno učenje
  - odločitvena drevesa (*Decision Trees*)
  - odločitvena pravila (*Decision Rules*)
  - učenje konceptov (*Concept Learning*)
  - indukcija logičnih programov (*ILP: Inductive Logic Programming*)
- Umetne nevronske mreže
  - večnivojske usmerjene NM
  - Kohonenove NM
  - Hopfieldove NM
  - Bayesove

Marko Bošker et al.

---

---

---

---

---

---

---

## Zahteve pri strojnem učenju

- Zanesljivost delovanja
  - velika klasifikacijska točnost
- Transparentnost naučenega znanja
  - eksplicitna simbolična predstavitev, razumljiva ekspertom
- Sposobnost pojasnjevanja
  - argumentiranje in podpora ekspertnim odločitvam
- Odpornost na "šum" v podatkih
  - delovanje ob manjkajočih, nepopolnih ali nenantančnih podatkih
  - problemi iz realnega sveta

Marko Bošker et al.

---

---

---

---

---

---

---

## Nekatera področja uporabe

- Medicina
  - Diagnostika in prognostika
- Industrija
  - Kontrola kvalitete
  - Procesna kontrola
- Upravljanje in odločanje
  - Analiza podatkov o poslovanju
- "Podatkovno rudarjenje" (*Data Mining*)

Marko Bošker et al.

---

---

---

---

---

---

---

## Viri



Tom M. Mitchell: *Machine Learning*.  
McGraw-Hill, 1997.



Igor Kononenko: *Strojno učenje, druga izdaja*.  
Založba FE in FRI, 2005.

Mario Boharec

---

---

---

---

---

---

---

---

## Kazalo

- Uvod
- Metode strojnega učenja
  - Atributno učenje
  - Učenje odločitvenih dreves
  - Statistične metode
  - Umetne nevronske mreže
  - HINT: učenje hierarhičnih modelov (DEX)

Mario Boharec

---

---

---

---

---

---

---

---

## Atributno učenje

### Podano:

- tabela "rešenih" primerov,
- opisanih z vrednostmi atributov  $A_1$  do  $A_N$
- in razredom  $C$

$C$	$A_1$	$A_2$	...	$A_N$
$C_1$	$V_{1,1}$	$V_{1,2}$	...	$V_{1,N}$
$C_2$	$V_{2,1}$	$V_{2,2}$	...	$V_{2,N}$
...				
$C_M$	$V_{M,1}$	$V_{M,2}$	...	$V_{M,N}$

**Naloga:** poiskati pravilo za razred  $C$  glede na vrednosti  $A_1$  do  $A_N$

Mario Boharec

---

---

---

---

---

---

---

---

## Odločitvena drevesa

- Vozlišča predstavljajo pogoje (teste)
  - praviloma: preverjanje vrednosti atributa
  - pogosto: binarna vejitev (pogoji tipa  $A = v$  ali  $A \geq v$ )
- Listi predstavljajo razrede
- Drevo omogoča klasifikacijo novih primerov:
  - začni pri korenju
  - potuj navzdol v skladu z rezultati testov
  - odgovor (razred) je v listu

Marko Bošker

---

---

---

---

---

---

## Učenje odločitvenih dreves

**TDIDT:** Top-Down Induction of Decision Trees

Oseba	Starost	Spol	Dohodki	Stranka
Ana Kranjc	32	Ž	10.000	da
Micka Kovač	53	Ž	1.000.000	da
Meta Novak	27	Ž	20.000	ne
Jana Bevc	55	Ž	20.000	da
Peter Dolenc	26	M	100.000	da
Janez Gorenc	50	M	200.000	da

primer iz portala Študenti Datalniki

Marko Bošker

---

---

---

---

---

---

## Klasifikacija in regresija

- **Klasifikacija:**
  - razred  $C$  je diskretna spremenljivka
  - pravilo (Oseba, Starost, Spol, Dohodki)  $\Rightarrow$  Stranka
- **Regresija:**
  - "razred"  $C$  je zvezna spremenljivka
  - pravilo (Oseba, Starost, Spol, Stranka)  $\Rightarrow$  Dohodki

Marko Bošker

---

---

---

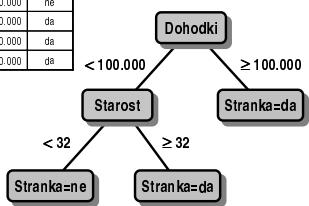
---

---

---

## Klasifikacijsko odločitveno drevo

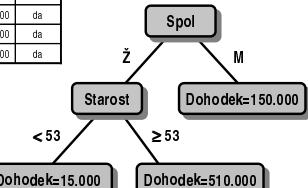
Oseba	Starost	Spol	Dohodki	Stranka
Ana Kranjc	32	Z	10.000	da
Micka Kovac	53	Z	1.000.000	da
Mela Novak	27	Z	20.000	ne
Jana Bevc	55	Z	20.000	da
Peter Dolenc	26	M	100.000	da
Janez Gorenc	50	M	200.000	da



Marko Bošnjak

## Regresijsko odločitveno drevo

Oseba	Starost	Spol	Dohodki	Stranka
Ana Kranjc	32	Z	10.000	da
Micka Kovac	53	Z	1.000.000	da
Mela Novak	27	Z	20.000	ne
Jana Bevc	55	Z	20.000	da
Peter Dolenc	26	M	100.000	da
Janez Gorenc	50	M	200.000	da



Marko Bošnjak

## Učenje odločitvenih dreves

### KLJUČNI KONCEPTI

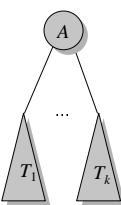
1. Gradnja drevesa
  - algoritem
  - izbiranje atributov
2. Preverjanje kakovosti drevesa
  - učna in testna množica
  - klasifikacijska točnost
3. Rezanje drevesa
  - rezanje naprej
  - rezanje nazaj

Marko Bošnjak

## Gradnja klasifikacijskega drevesa

### ALGORITEM

- Če vsi učni primeri pripadajo istemu razredu  $C$ , potem je rezultat list  $C$
- Sicer
  - Izberi najboljši atribut  $A$  (ali najboljšo delitev po  $A$ )
  - Razdeli učno množico glede na vrednosti  $A$
  - Rekurzivno zgradi poddrevesa  $T_1..T_k$  za vsako podmnožico
  - Rezultat je drevo z vozliščem  $A$  in poddrevesi  $T_1..T_k$



izvorno iz: Bojan Čestnik: Strojno učenje

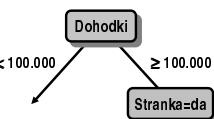
Marko Batagelj

## Primer

Oseba	Starost	Spol	Dohodki	Stranka
Ana Kranjc	32	Z	10.000	da
Mirka Kovac	53	Z	1.000.000	da
Mira Novak	27	Z	20.000	ne
Jana Bevc	55	Z	20.000	da
Peter Dolenc	26	M	100.000	da
Janez Gorenc	50	M	200.000	da

Vsi primeri vistem razredu?

Najboljša delitev?



Oseba	Starost	Spol	Dohodki	Stranka
Ana Kranjc	32	Z	10.000	da
Mirka Kovac	53	Z	1.000.000	da
Peter Dolenc	26	M	100.000	da

Oseba	Starost	Spol	Dohodki	Stranka
Mira Novak	27	Z	20.000	ne
Jana Bevc	55	Z	20.000	da

Marko Batagelj

## Izbiranje atributov (delitev)

**ZAHTEVA:** Delitev na čim bolj "čiste" podmnožice

### MERE "NEČISTOČE"

za dva razreda,  $p(C_1)=p_1, p(C_2)=p_2$

- Entropija  $E$ :  $-p_1 \log_2 p_1 - p_2 \log_2 p_2$
- Napaka prevladujočega razreda:  $1 - \max(p_1, p_2)$
- Indeks Gini:  $1 - (p_1^2 + p_2^2)$

**INFORMACIJSKI PRISPEVEK:** Koliko pridobimo ob delitvi?

- $\text{Gain}(S, A) = E(S) - \sum_v |S_v| / |S| E(S_v)$
- maksimiziramo Gain

izvorno iz: Bojan Čestnik: Strojno učenje

Marko Batagelj

## Primer

(1 od 3)

Vreme	Temp	Vlaga	Veter	Tenis
sončno	vruče	visoka	ne	ne
sončno	vruče	visoka	da	ne
oblačno	vruče	visoka	ne	da
dež	zmerino	visoka	ne	da
dež	hladno	norm	ne	da
dež	hladno	norm	da	ne
oblačno	hladno	norm	da	da
sončno	zmerino	visoka	ne	ne
sončno	hladno	norm	ne	da
dež	zmerino	norm	ne	da
sončno	zmerino	norm	da	da
oblačno	zmerino	visoka	da	da
oblačno	vruče	norm	ne	da
dež	zmerino	visoka	da	ne

Vreme? sončno [2+, 3-]  $E=0,97$   
 oblačno [4+, 0-]  $E=0$   
 dež [3+, 2-]  $E=0,97$

Vlaga? visoka [3+, 4-]  $E=0,99$   
 norm [6+, 1-]  $E=0,59$

Veter? ne [6+, 2-]  $E=0,81$   
 da [3+, 3-]  $E=1,00$

Marko Boškarec

## Primer

(2 od 3)

Veter? ne [6+, 2-]  $E=0,81$   
 da [3+, 3-]  $E=1,00$

$$\text{Gain}(S, A) = E(S) - \sum_v |S_v| / |S| \cdot E(S_v)$$

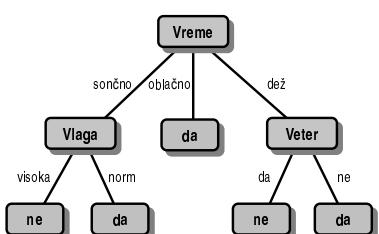
$$E(S) = E(9+, 5-) = -(9/14)\log_2(9/14) - (5/14)\log_2(5/14) = 0,94 \\ E(S_{\text{Veter}=ne}) = 0,81 \\ E(S_{\text{Veter}=da}) = 1,00 \\ \text{Gain}(S, \text{Veter}) = 0,94 - (8/14)0,81 - (6/14)1,00 = \mathbf{0,048}$$

$$\begin{aligned} \text{Gain}(S, \text{Vreme}) &= \mathbf{0,246} \text{ (max)} \\ \text{Gain}(S, \text{Vlaga}) &= \mathbf{0,151} \\ \text{Gain}(S, \text{Temp}) &= \mathbf{0,029} \end{aligned}$$

Marko Boškarec

## Primer

(3 od 3)



Marko Boškarec

## Mere kvalitete odločitvenih dreves

### Klasifikacijska točnost:

Kako točno drevo klasificira nove primere?  
Kakšna je točnost v primerjavi z *apriorno* ("naivni klasifikator")?

### Razumljivost:

Ali ekspert razume drevo in njegovo vsebino?  
Ali ga lahko interpretira, utemelji?

### Velikost:

Povezano z razumljivostjo: zaželena čim manjša drevesa!

Marko Bošker

---

---

---

---

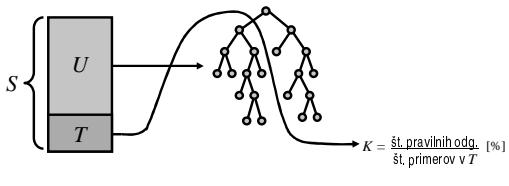
---

---

## Klasifikacijska točnost

### Postopek gradnje in preverjanja klasifikatorja:

- množico primerov  $S$  razdelimo na:
  - **učno** množico  $U$  (npr. 70%) in
  - **testno** množico  $T$  (30%)
- zgradimo drevo upoštevajoč samo  $U$
- na  $T$  preverimo **točnost klasifikacije** = delež pravilnih klasifikacij



Marko Bošker

---

---

---

---

---

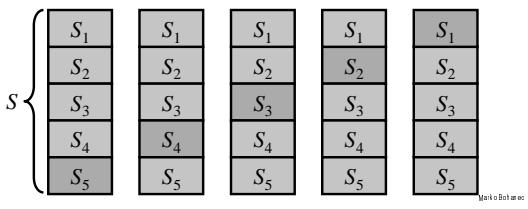
---

## Prečno preverjanje

Slov: Prečno preverjanje, npr. "10-kratno prečno preverjanje"  
Angl: "Cross validation", e.g., "10-fold cross validation"

### Postopek preverjanja klasifikatorja:

- celoto množico primerov  $S$  zaključno razdelimo na  $n$  čim bolj enakih delov
- $n$  krat ponovimo: učimo na  $n-1$  delih, testiramo na preostalem enem delu
- celotna točnost klasifikacije je povprečje  $n$  takšnih klasifikacij



Marko Bošker

---

---

---

---

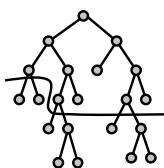
---

---

## Rezanje dreves

Spodnji deli drevesa (okrog listov) so manj zanesljivi:

- manjše število učnih primerov
- prevelika prilagoditev podatkom (*overfitting*)



Rezanje (*pruning*):

- Naprej: predčasno ustavimo gradnjo
- Nazaj: drevo zgradimo do konca, nato rezemo manj zanesljive dele (boljje)

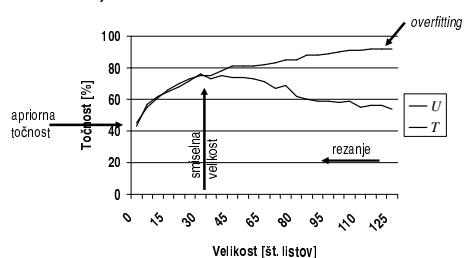
Pridobimo:

- Manjše drevo – večja preglednost in razumljivost
- Večja točnost na testni množici primerov

Marko Bošker et al.

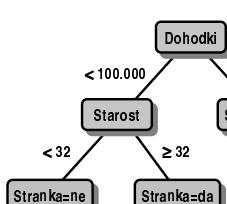
## Rezanje dreves

Klasifikacijska točnost v odvisnosti od velikosti drevesa



Marko Bošker et al.

## Od dreves k pravilom



**PRAVILA**

```

if Dohodki ≥ 100.000
then Stranka=da
if Dohodki < 100.000 and
Starost ≥ 32
then Stranka=da
if Dohodki < 100.000 and
Starost < 32
then Stranka=ne
  
```

**ODLOČITVENI SEZNAM**

```

if Dohodki ≥ 100.000
then Stranka=da
else if Starost ≥ 32
then Stranka=da
else Stranka=ne
  
```

Marko Bošker et al.

## Nekaj učnih algoritmov in programov

### KLASIFIKACIJA

- ID3 (Quinlan 1979)
- CART (Breiman et al. 1984)
- Assistant (Cestnik et al. 1987)
- C4.5 (Quinlan 1993)
- C5.0, See5 (RuleQuest)
  - Weka (Waikato University, NZ)
  - Orange (Univerza v Ljubljani, FRI)
- Cubist (RuleQuest)

### REGRESIJA

Marko Bošker

---



---



---



---



---



---



---



---



---

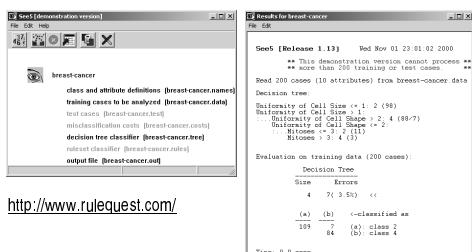


---



---

## See5 (RuleQuest)



Marko Bošker

---



---



---



---



---



---



---



---



---



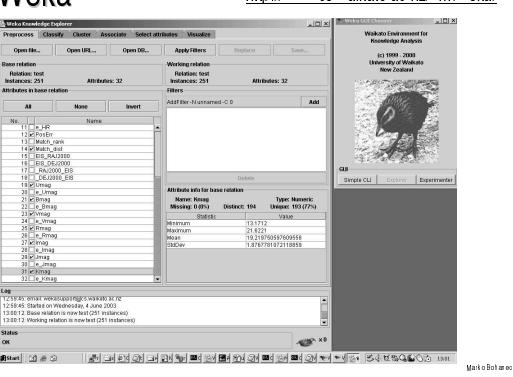
---



---

## Weka

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




---



---



---



---



---



---



---



---



---



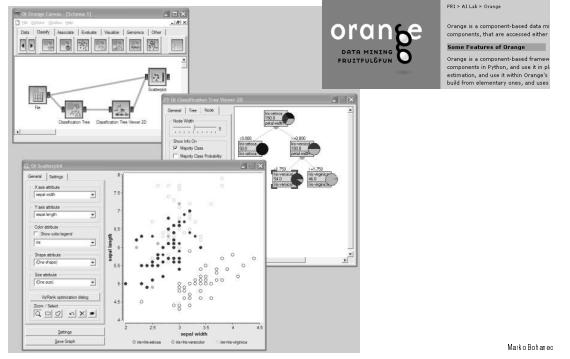
---



---

## Orange

<http://magix.fri.uni-lj.si/orange/>



## Kazalo

- Uvod
- Metode strojnega učenja
  - Atributno učenje
  - Učenje odločitvenih dreves
  - Statistične metode
    - Bayesovo pravilo
    - k-NN: k najbližjih sosedov
  - Umetne nevronske mreže
  - HINT: učenje hierarhičnih modelov (DEX)

## Bayesovo pravilo

- *Bayesovo pravilo:* verjetnost razreda  $C$  pri pogoju  $V_1..V_n$

$$P(C|V_1..V_n) = P(C) \frac{P(V_1..V_n|C)}{P(V_1..V_n)}$$

- “*Naivno*” *Bayesovo pravilo:* predpostavka neodvisnosti  $V_1..V_n$

$$P(C|V_1..V_n) = P(C) \prod_{i=1}^n \frac{P(C|V_i)}{P(C)}$$

## Primer

Vreme	Temp	Višaga	Veter	Tenis
sončno	vroče	visoka	ne	ne
sončno	vroče	visoka	da	ne
oblačno	vroče	visoka	ne	da
dež	zmereno	visoka	ne	da
dež	hladno	norm	ne	da
dež	hladno	norm	da	ne
oblačno	hladno	norm	da	da
sončno	zmereno	visoka	ne	ne
sončno	hladno	norm	ne	da
dež	zmereno	norm	ne	da
sončno	zmereno	norm	da	da
oblačno	zmereno	visoka	da	da
oblačno	vroče	norm	ne	da
dež	zmereno	visoka	da	ne

Vreme	Temp	Višaga	Veter	Tenis	?
sončno	hladno	norm	da	?	

$$P(C|V_1..V_n) = P(C) \prod_{i=1}^n \frac{P(C|V_i)}{P(C)}$$

$$P(da|sončna norm da) = \\ = P(da) \cdot \frac{P(da|sončna) P(da|norm) P(da|norm)}{P(da)} \cdot \frac{P(da)}{P(da)}$$

$$= \frac{9}{14} \cdot \frac{2}{5} \cdot \frac{14}{9} \cdot \frac{3}{4} \cdot \frac{14}{9} \cdot \frac{6}{7} \cdot \frac{14}{9} = 0,48 \\ \Rightarrow \text{Tenis=da}$$

$$P(ne|sončna norm da) = \\ = P(ne) \cdot \frac{P(ne|sončna) P(ne|norm) P(ne|norm)}{P(ne)} \cdot \frac{P(ne)}{P(ne)}$$

$$= \frac{5}{14} \cdot \frac{3}{5} \cdot \frac{14}{9} \cdot \frac{4}{5} \cdot \frac{14}{7} \cdot \frac{3}{6} = 0,24$$

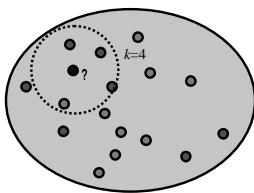
Marko Boškar ec

## K-najbližjih sosedov ( $k$ -NN)

- Pri klasifikaciji primera  $X = V_1, V_n$  v učni množici poiščemo *k najbližjih sosedov*
- Primer  $X$  klasificiramo v večinski razred  $C$

Definirati je potrebno:

- $k$
- mero razdalje med primeri



Marko Boškar ec

## Primer

Vreme	Temp	Višaga	Veter	Tenis	Razd
sončno	vroče	visoka	ne	ne	4
sončno	vroče	visoka	da	ne	3
oblačno	vroče	visoka	ne	da	5
dež	zmereno	visoka	ne	da	5
dež	hladno	norm	ne	da	3
dež	hladno	norm	da	ne	2
oblačno	hladno	norm	da	da	1
sončno	zmereno	visoka	ne	ne	3
sončno	hladno	norm	ne	da	1
dež	zmereno	norm	ne	da	4
sončno	zmereno	norm	da	da	1
oblačno	zmereno	visoka	da	da	3
oblačno	vroče	norm	ne	da	4
dež	zmereno	visoka	da	ne	4

Vreme	Temp	Višaga	Veter	Tenis	?
sončno	hladno	norm	da	?	

$k = 4$ , razdalja "Manhattan"

$\Rightarrow \text{Tenis=da (3:1)}$

Marko Boškar ec

## Kazalo

- Uvod
- Metode strojnega učenja
  - Atributno učenje
  - Učenje odločitvenih dreves
  - Statistične metode
    - Bayesovo pravilo
    - k-NN: k najbližjih sosedov
  - Umetne nevronske mreže
  - HINT: učenje hierarhičnih modelov (DEX)

Marko Bošker

---



---



---



---



---



---

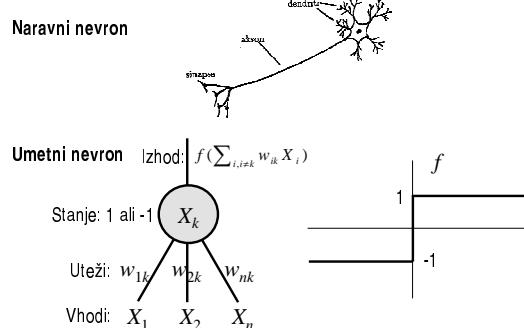


---



---

## Umetne nevronske mreže




---



---



---



---



---



---

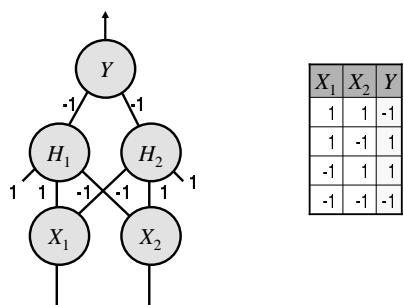


---



---

## Primer: XOR (ekskluzivni "ali")




---



---



---



---



---



---



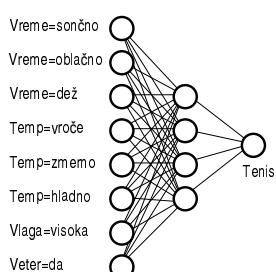
---



---

## Primer

Vreme	Temp	Vlag a	Veter	Tenis
sončno	vroče	visoka	ne	ne
sončno	vroče	visoka	da	ne
oblačno	vroče	visoka	ne	da
dež	zmerino	visoka	ne	da
dež	hladno	norm	ne	da
dež	hladno	norm	da	ne
oblačno	hladno	norm	da	da
sončno	zmerino	visoka	ne	ne
sončno	hladno	norm	ne	da
dež	zmerino	norm	ne	da
sončno	zmerino	norm	da	da
oblačno	zmerino	visoka	da	da
oblačno	vroče	norm	ne	da
dež	zmerino	visoka	da	ne



Marko Bošker et al.

## Lastnosti umetnih nevronskeih mrež

- Matematična osnova  
linearna algebra
- Paralelizem  
nevroni hkrati izvajajo operacije
- Večsmerno izvajanje  
vhod je lahko tudi izhod in obratno
- Robustnost  
majhna občutljivost na okvare nevronov in sinaps
- Učenje  
s spontanim spreminjanjem uteži  
razmeroma počasno
- Razlaga odločitve  
slaba, praktično nemogoča

Marko Bošker et al.

## Kazalo

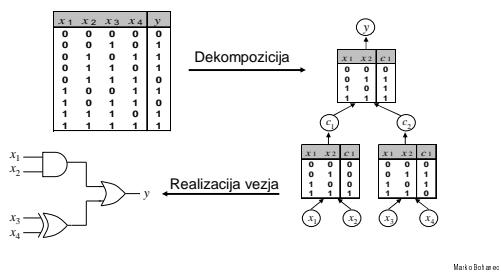
- Uvod
- Metode strojnega učenja
  - Atributno učenje
  - Učenje odločitvenih dreves
  - Statistične metode
    - Bayesovo pravilo
    - k-NN: k najbližjih sosedov
  - Umetne nevronske mreže
  - HINT: učenje hierarhičnih modelov (DEX)

Marko Bošker et al.

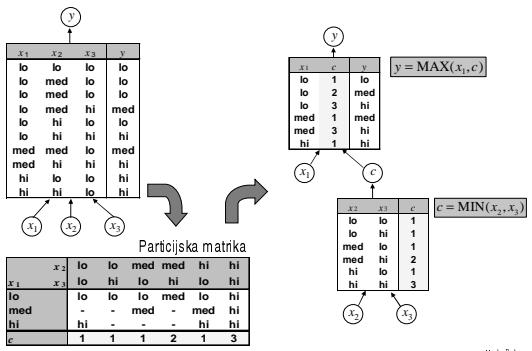
## HINT: Učenje hierarhičnih modelov

Cilj: učenje hierarhičnih modelov tipa DEX

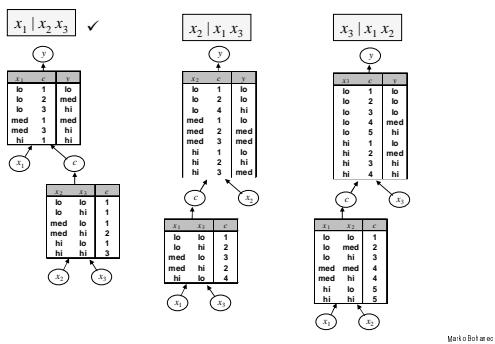
Algoritam temelji na Ashenhurst-Curtisovi dekompoziciji logičnih funkcij



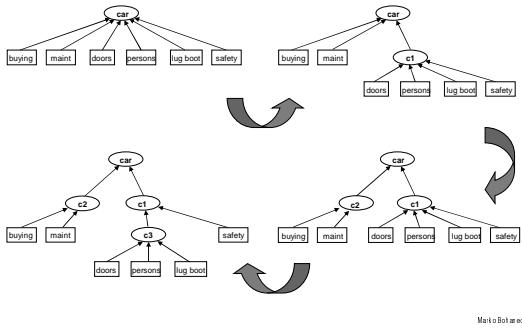
## HINT: Osnovni korak



## HINT: Izbira najboljše delitve



## HINT: Oblikovanje hierarhije konceptov



---

---

---

---

---

---

## Strojno učenje: Povzetek

- Rešujejo problem "ozkega grla" (*Feigenbaum's Bottleneck*)
- Učenje je "avtomatizirano" (skupaj z ekspertom), ne "avtomatsko"
- Številne metode: statistične, simbolične, nevronske mreže
- Pomembni koncepti:
  - klasifikacija, regresija
  - klasifikacijska točnost
  - razumljivost, velikost modelov
  - učna in testna množica
  - prilagoditev podatkom (*overfitting*)
- Odločitvena drevesa: razširjena, pogosta uporaba; koncept rezanja
- Statistične metode: dobre, a ne gradijo opisov konceptov
- Nevronske mreže: zanimive lastnosti, zahtevno učenje, problem razumljivosti
- HINT: učenje konceptov po načelu razgradnje funkcij

Mark o Bot ar ec

---

---

---

---

---

---