

# 1 Osnove

## 1.1 Kaj je umetna inteligenca?

- **cilji:** Razumeti in zgraditi inteligentne sisteme na osnovi razumevanja človeškega razmišljanja, sklepanja, učenja in komuniciranja.

# 2 Strojno Učenje

## 2.1 Kaj je strojno učenje?

Je področje umetne inteligence, ki raziskuje kako se lahko algoritmi samodejno izboljšujejo ob pridobivanju izkušenj.

## 2.2 Vrste učenja:

- **Nadzorovano učenje** *supervised learning*: Učni primeri so označeni in podani kot vrednosti vhodov in izhodov. Učimo se funkcije, ki vhode preslika v izhode. (npr. odločitveno drevo)
- **Nenadzorovano učenje** *unsupervised learning*: Učni primeri niso označeni → nimajo ciljne spremenljivke. Učimo se iz vzorcev v podatkih. (npr. gručenje)
- **Spodbujevalno učenje** *reinforcement learning*: Inteligentni agent se uči iz zaporedja nagrad in kazni.

## 2.3 Nadzorovano učenje

Podano imamo množico **učnih** primerov:

$$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$$

kjer je vsak  $y_j$  vrednost neznane funkcije  $y = f(x)$ . Nasa naloga je posikati hipotetično funkcijo  $h$ , ki je najboljši možen približek funkciji  $f$ .

Locimo dve vrsti problemov:

1. **Klasifikacijski:**  $y_j$  je *diskretna(kategoricna)* spremenljivka
  - $y$  pripada **koncnemu naboru vrednosti** (diskretna spremenljivka)
  - $y$  imenujemo **razred** (class)
2. **Regresijski:**  $y_j$  je *zvezna* spremenljivka
  - $y$  je stevilo (običajno  $y \in R$ , je zvezna spremenljivka)
  - $y$  imenujemo **oznacza** (label)

### 2.3.1 Prostor in evalviranje hipotez

Denimo da imamo:

- binarno klasifikacijo
- $n$  binarnih atributov

Iz tega sledi:

- $2^n$  različnih učnih primerov
- $2^{2^n}$  hipotez (celotno odločitveno drevo)

Pomembni kriteriji pri *evalviranju* hipotez:

- **konsistentnost** hipotez s (učnimi) primeri
- **splosnost** točnost za nevidene primere

- **razumljivost** (*interpretability, comprehensibility*) hipotez

Poznamo 4 razrede za ocenjevanje uspešnosti pri klasifikaciji na podlagi njihove **točnosti**:

- **TP** - pravilno pozitivno klasificirani primeri
- **TN** - pravilno negativno klasificirani primeri
- **FP** - napacno pozitivno klasificirani primeri
- **FN** - napacno negativno klasificirani primeri

**Klasifikacijska točnost** je potem definirana:

$$CA = \frac{TP+TN}{TP+TN+FP+FN} = \frac{TP+TN}{N}$$

Poznamo dva tipa atributov:

### 1. diskretni atributi:

- **nominalni** npr. ['sončno', 'dezeveno']
- **ordinalni** npr. ['nizko', 'srednje', 'visoko']

Odločitvena drevesa delijo prvotno učno množico na vse manjše podmnožice

### 2. zvezni atributi:

Delitev podmnožice glede na smislno mejo izbranega atributa

### 2.3.2 Odločitveno drevo

Ponazarja relacijo med vhodnimi *vrednostmi/atributi* in *odločitvojo/ciljno* spremenljivko.

Z **notranjimi vozlišci** opravljamo test glede na vrednost posameznega atributa. Na koncu pridemo do **lista**, ki nam s poroči odločitev (vrednost ciljne spremenljivke). Konjunkcijo pogojev v *notranjih vozlišcih* katera vodi do *lista* imenujemo **pot**.

Gradnja odločitvenega drevesa: Nas cilj je zgraditi **cim manjše drevo**, ki je **konsistentno** z učnimi podatki.

**Hevristični pozrezní algoritem - TDIDT** s strategijo **razveji in omeji**:

- Izberi najbolj pomemben atribut - tisti, ki najbolj odločilno vpliva na klasifikacijo primera in razdeli primere v poddrevesa glede na njegove vrednosti
- rekurzivno ponovi za vsa drevesa
- ce vsi elementi v listu pripadajo istemu razredu ali vozlišča ni možno deliti naprej(ni razpoložljivih atributov), ustavi gradnjo

**Kratovidnost TDIDT:** Ker je TDIDT pozrezní algoritem, ki "lokalno" izbira najboljši atribut, ne upsteva kako dobro drugi algoritmi doplnjujejo izbrani atribut.

### 2.3.2.1 Izbor najbolj pomembnega atributa in informacijski prispevek

Najboljši atribut je tisti, ki razdeli učno množico v najbolj čiste podmnožice. Uporabimo lahko mero entropije:

$$H(X) = \sum_{i=1}^n p_i I_i = - \sum_{i=1}^n p_i \log_2 p_i$$

Zanima nas **znízanje** entropije ( *nedolocenosti* ) ob delitvi učne množice glede na vrednosti atributa  $A$ .

Definirajmo **informacijski prispevek** na takšen način, da najbolj informativni atribut **maksimizira informacijski prispevek** oz. minimizira  $I_{res}$ .

$$Gain(A) = I - I_{res}(A) \\ I_{res} = - \sum_{v_i \in A} p_{v_i} \sum_c p(c|v_i) \log_2 p(c|v_i)$$

2.3.2.2 Vecvrednostni atributi

Tezava z atributi, ki imajo vec kot dve vrednosti: Informacijski prispevek precenjuje njihovo kakovost(entropija je visja na racun vecjega stevila vrednosti in ne na racun kakovosti atributa) resitve:

- normalizacija informacijskega prispevka: **relativni informacijski prispevek** ali IGR (information gain ratio)

$$Gain(a) = I - I_{res}(A), I(A) = - \sum_v p_v \log_2 p_v$$
$$GainRatio(A) = \frac{Gain(A)}{I(A)} = \frac{I - I_{res}(A)}{I(A)} \text{ Oba zelimo maksimizirati}$$

- uporaba **alternativnih mer**: npr. **Gini index** *Ocena pricakovane klasifikacijske napake, vsota produktov verjetnosti razredov*

$$Gini = \sum_{c_1 \neq c_2} p(c_1)p(c_2)$$
$$Gini(A) = \sum_v p(v) \sum_{c_1 \neq c_2} p(c_1|v)p(c_2|v)$$

- **binarizacija** atributov: Je alternativa za reševanje problematike z vecvrednostmi atributi. Prednosti binarizacije so manjša vejanja drevesa, kar je statisticno bolj zanesljivo. Razlicni nacini binarizacije atributa lahko nastopajo kot samostojni atributi, ki se v drevesu pojavijo veckrat.

2.3.3 Privzeta tocnost in Pristranost na ucnii mnozici

Smiselna mera za **Privzeto tocnost** odlocitvenega drevesa je **verjetnost vecinskega razreda** v ucnii mnozici. Drevo je uporabno,

ce je njegova tocnost **visja** od privzete tocnosti.  
npr. [#Da, #Ne] = [3, 7] → verjetnost vecinskega razreda: 7/10  
Nas cilj je maksimizirati pricakovano tocnost na testnih podatkih vendar se zelimo izogniti pretiranemu prilaganju. Zato obicajno podatke razdelimo na **ucno** (70%) in **testno** mnozico (30%).

2.3.3.1 Ucenje dreves iz sumnih podatkov

V primeru da podatki niso popolni(premalo primerov / atributov) ali napake se lahko pojavijo tezave:

- **Ucenje suma** in ne dejanske aproksimacijske funkcije
- Pretirano prilaganje vodi v **prevelika drevesa** → overfitting
- **Slaba razumljivost** dreves

Posledica: **nizja klasifikacijska tocnost** na novih nevidenih podatkih.

2.3.3.2 Rezanje odlocitvenih dreves

Resujemo problem prevelikega prilaganja ucnim podatkom. Nizji deli drevesa predstavljajo vecje lokalno prilaganje ucnim podatkom, ki so lahko posledica suma. Poznamo dve metodi rezanja dreves:

1. Rezanje z zmanjševanjem napake (**REP**)
2. Rezanje z minimizacijo napake (**MEP**)