

1 Osnove

1.1 Kaj je umetna inteligenca?

- **cilji:** Razumeti in zgraditi inteligentne sisteme na osnovi razumevanja cloveskega *razmisljanja, sklepanja, ucenja* in *kommuniciranja*.

2 Strojno Učenje

2.1 Kaj je strojno ucenje?

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

2.2 Vrste ucenja:

- **Nadzorovano ucenje** *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 ucenje** *unsupervised learning*: Učni primeri niso označeni → nimajo ciljne spremenljivke. Učimo se iz vzorcev v podatkih. (npr. gručenje)
- **Spodbujevalno ucenje** *reinforcement learning*: Inteligentni agen se uči iz zaporedja nagrad in kazni.

2.3 Nadzorovano ucenje

Podano imamo množico **ucnih** 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 (obicačno $y \in R$, 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 ucnih primerov
- 2^{2^n} hipotez (celotno odločitveno drevo)

Pomembni kriteriji pri *evalviranju* hipotez:

- **konsistentnost** hipotez s (ucnimi) 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. [’soncno’, ’dezevno’]
- **ordinalni** npr. [’nizko’, ’srednje’, ’visoko’]

Odločitvena drevesa delijo prvotno ucno 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 ucnimi podatki.

Hevristicni pozreznis 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 vozlišca ni možno deliti naprej (ni razpoložljivih atributov), ustavi gradnjo

Kratovidnost TDIDT: Ker je TDIDT pozreznis 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 ucno množico v najbolj ciste 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 **znizanje** entropije (*nedolocnosti*) ob delitvi ucne 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 Vcevrednostni atributi

Tezava z atributi, ki imajo več kot dve vrednosti: Informacijski prispevek preceňuje njihovo kakovost (entropija je visja na račun večjega števila vrednosti in ne na račun kakovosti atributa)

resitve:

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

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

Oba želimo maksimizirati

- uporaba **alternativnih mer**: npr. **Gini index** *Ocena pričakovane 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 vcevrednostnimi atributi. Prednosti binarizacije so manjša vejanja drevesa, kar je statistično bolj zanesljivo. Različni načini binarizacije atributa lahko nastopajo kot samostojni atributi, ki se v drevesu pojavijo večkrat.

2.3.2.3 Privzeta točnost in Pristranost na učni množici

Smiselna mera za **Privzeto točnost** odločitvenega drevesa je **verjetnost večinskega razreda** v učni množici. Drevo je uporabno, če je njegova točnost **visja** od privzete točnosti.

npr. $[\#Da, \#Ne] = [3, 7] \rightarrow$ verjetnost večinskega razreda: $7/10$

Nas cilj je maksimizirati pričakovano točnost na testnih podatkih vendar se želimo izogniti pretiranemu prilaganju. Zato običajno podatke razdelimo na **učno** (70%) in **testno** množico (30%).

2.3.2.4 Učenje dreves iz sumnih podatkov

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

- **Učenje suma** in ne dejanske aproksimacijske funkcije
- Pretirano prilaganje vodi v **prevelika drevesa** \rightarrow overfitting
- **Slaba razumljivost** dreves

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

2.3.2.5 Rezanje odločitvenih dreves

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

1. **Rezanje vnaprej**: uporaba dodatnega kriterija za zaustavitev gradnje drevesa glede na obseg suma. Je hitrejši pristop, vendar kratkoviden (uposteva samo zgornji del drevesa)
2. **Rezanje nazaj**: Rezanje, ki po gradnji celotnega drevesa, odstrani manj zanesljive dele drevesa (opisujejo sum, zgrajeni iz manj podatkov in z manj informirani atributi) Je počasnejše, in uposteva informacijo iz **celega drevesa**.

1. Rezanje z zmanjševanjem napake (REP)

Uporablja posebno reazalno (validacijsko) množico:

- **Učna množica** (70%):
 - množica za gradnjo *growing set* (70%)
 - reazalna množica *pruning set* (30%)
- **Testna množica** (30%)

Postopek:

- (a) Potuj od listov navzgor (pricni s starsi listov)
- (b) Za vsako vozlišče **v** izračunaj **dobitek rezanja**
Dobitek rezanja izračunamo: st. napacnih klasifikacij v drevesu **T** - st. napacnih klasifikacij v vozlišču **v**.
- (c) Če je dobitek ≥ 0 , obrezi in nadaljuj postopek s starsim, sicer ustavi postopek

2. Rezanje z minimizacijo napake (MEP)

Uporablja množico za gradnjo drevesa in ne ločene reazalne množice. Postopek: Za vozlišče **v** izračunamo:

- (a) **staticno napako** (verjetnost klasifikacije v napacnem razred)

$$e(v) = p(\text{razred} \neq C|v), C \text{ je večinski razred v } v$$

- (b) **vzratno napako** (*backed-up error*)

$$\sum_i p_i E(T_i) = p_1 E(T_1) + p_2 E(T_2) + \dots$$

Rezemo, če je **staticna napaka** $<$ **vzratna napaka**.

Napaka **optimalno** obrezanega drevesa je torej:

$$E(T) = \min(e(v), \sum_i p_i E(T_i)) \text{ in } E(T) = e(v), \text{ ce je } v \text{ list.}$$

Namesto **minimizacije** napake E, lahko problem obrnemo in **maksimiziramo** točnost CA.

2.3.2.6 Ocenjevanje verjetnosti

1. Laplaceova ocena verjetnosti

$$p = \frac{n+1}{N+k}$$

n - st. primerov, ki pripadajo razredu C

N - st. vseh primerov

k - st. vseh razredov

k je problematičen parameter, saj ocena ne upošteva apriorne verjetnosti

2. m-ocena verjetnosti

$$p = \frac{n+p_a m}{N+m} = p_a \frac{m}{N+m} + \frac{n}{N} \frac{N}{N+m}$$

p_a - Apriorna verjetnost razreda C

m - parameter ocene (vpliva na delež upoštevanja apriorne verjetnosti)

Ce imamo malo suma, potem m nastavimo na majno vrednost in imamo malo rezanja. Obratno v primeru ce imamo veliko suma. Gre se za posplošitev Laplaceove ocene za $m = k$ in $p_a = 1/k$.