

OSNOVE UMETNE INTELIIGENCE

2018/19

regresija
linearne in lokalne metode
ocenjevanje učenja

Pregled



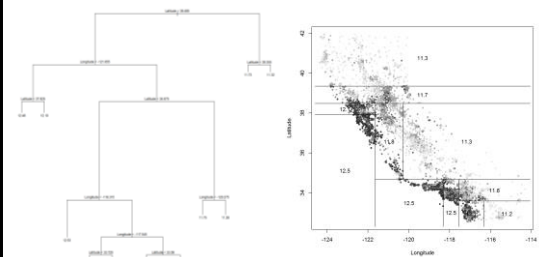
- strojno učenje
 - uvod v strojno učenje
 - učenje odločitvenih dreves
 - učenje dreves iz šumnih podatkov (rezanje dreves)
 - regresijska drevesa
 - linearni modeli
 - metoda k najbližjih sosedov
 - ocenjevanje učenja

Regresijska drevesa

- zvezna ciljna spremenljivka – regresijski problem
- regresijska drevesa so podobna odločitvenim drevesom, le za regresijske probleme
- sistemi: CART (Breiman et al. 1984), RETIS (Karalič 1992), M5 (Quinlan 1993), WEKA (Witten and Frank, 2000)
- listi v regresijskem drevesu predstavljajo:
 - predstavljajo povprečno vrednost označb ("razreda") primerov v listu
 - preprost napovedni model (npr. linearna regresija) za nove primere



Regresijska drevesa



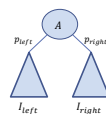
Gradnja regresijskih dreves

- atribut delimo glede na izbrano mejno vrednost
- drugačna mera za merjenje nedoločenosti/nečistoče: srednja kvadratna napaka v vozišču v:

$$MSE(v) = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

- cilj: minimiziramo rezidualno nedoločenost po delitvi primerov glede na vrednosti atributa A
- pričakovana rezidualna nečistost

$$I_{res}(A) = p_{left} \cdot I_{left} + p_{right} \cdot I_{right}$$



Pregled



- strojno učenje
 - uvod v strojno učenje
 - učenje odločitvenih dreves
 - učenje dreves iz šumnih podatkov (rezanje dreves)
 - regresijska drevesa
 - linearni modeli
 - metoda k najbližjih sosedov
 - ocenjevanje učenja

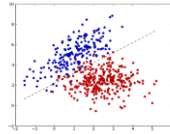
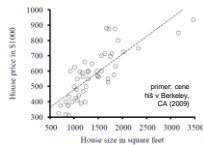
Linearni modeli

- uporaba pri **klasifikaciji** (kot separator razredov) in **regresiji** (kot prilaganje skozi podane točke)
- linearni model z **eno odvisno** spremenljivko (angl. *univariate linear model*):

$$h(x) = w_1 x + w_0$$

w_0 in w_1 sta **uteži** (angl. *weights*) spremenljivk (koeficienta)

- linearna regresija**: postopek iskanja funkcije $h(x)$ (oziroma uteži w_0 in w_1), ki se najbolj prilaga učnim podatkom



Linearna regresija

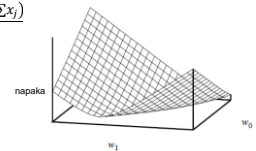
- optimizacijo izvedemo z minimizacijo srednje kvadratne napake:

$$\text{napaka}(h) = \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2$$

- prostor koeficientov je konvexen, lokalni minimumi ne obstajajo (samo globalni)
- obstaja analitična rešitev:

$$w_1 = \frac{N(\sum x_j y_j) - (\sum x_j)(\sum y_j)}{N(\sum x_j^2) - (\sum x_j)^2}$$

$$w_0 = \frac{\sum y_j - w_1(\sum x_j)}{N}$$



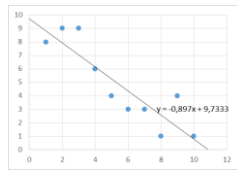
Linearna regresija

- primer linearne regresije

x_j	y_j	$x_j y_j$	x_j^2
1	8	8	1
2	9	18	4
3	9	27	9
4	6	24	16
5	4	20	25
6	3	18	36
7	3	21	49
8	1	8	64
9	4	36	81
10	1	10	100
$\sum x_j = 55$	$\sum y_j = 48$	$\sum x_j y_j = 190$	$\sum x_j^2 = 385$

$$w_1 = \frac{N(\sum x_j y_j) - (\sum x_j)(\sum y_j)}{N(\sum x_j^2) - (\sum x_j)^2} = \frac{10 \cdot 190 - 55 \cdot 48}{10 \cdot 385 - 55^2} = -0.897$$

$$w_0 = \frac{\sum y_j - w_1(\sum x_j)}{N} = \frac{48 - (-0.897) \cdot 55}{10} = 9.733$$



Posplošitev v več dimenzij

- možna je posplošitev v **višje število dimenzij** – več neodvisnih spremenljivk (atributov) (angl. *multivariate linear regression*)

$$h(x) = w_0 + \sum_i w_i x_{i,1}$$

kjer so w_i uteži (koeficienti), $x_{i,1}$ pa i -ta spremenljivka (atribut) primera x_j

- uteži lahko določimo **analitično**: $w = (X^T X)^{-1} X^T y$

kjer je X matrika s podatki (vrstice – učni primeri, stolpci – atributi), y pa vektor z vrednostmi odvisnih spremenljivk primerov

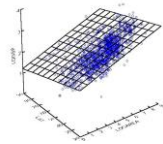
- v praksi se odločamo za iskanje koeficientov z **gradientnim spustom**

$w \leftarrow$ naključna začetna rešitev
ponavljaj do konvergence

za vsak w_i v w :

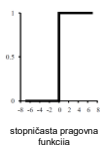
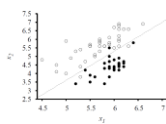
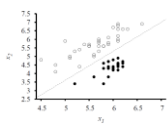
$$w_i \leftarrow w_i - \alpha \frac{\partial}{\partial w_i} \text{napaka}(w)$$

- problem s pretiranim prilagajanjem, regularizacija



Linearni modeli pri klasifikaciji

- linearni model se uporablja za ločevanje primerov, ki pripadajo različnim razredom
- iščemo **odločitveno mejo** (angl. *decision boundary*) oz. **linearni separator** (obstaja samo pri linearno ločljivih problemih)
- za spodnji primer je linearno separator lahko funkcija $-4.9 + 1.7x_1 - x_2 = 0$
- hipoteza je torej: $h(x) = \text{prag}(w \cdot x)$, kjer $\text{prag}(z) = \begin{cases} 1 & z \geq 0 \\ 0 & \text{sicer} \end{cases}$



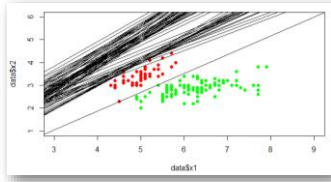
Linearni modeli pri klasifikaciji

- možnih ustreznih premic je več
- preprosto iskanje rešitve – **stohastični gradientni spust** s posodabljanjem uteži
- za vsak učni primer (x, y) izvedi posodobitev uteži:

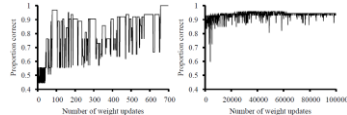
$$w_i \leftarrow w_i + \alpha (y - h(x)) \times x_i$$
 kjer so w_i uteži (koeficienti), α pa vpliva na hitrost spremembe (korak)
- intuicija:
 - če $y = h(x)$, potem se w_i ne spremeni
 - če $y = 1$ in $h(x) = 0$ (**prenizka** vrednost hipoteze), potem se za pozitiven x_i utež **poveča** in za negativen x_i utež **zmanjša**
 - če $y = 0$ in $h(x) = 1$ (**previsoka** vrednost hipoteze), potem se za pozitiven x_i utež **zmanjša** in za negativen x_i utež **poveča**
- algoritem lahko pri ustreznem α najde optimalno rešitev tudi za linearno neločljive podatke
- smiselna izbirljavanja: logistična pragovna funkcija

Linearni modeli pri klasifikaciji

- demo



konvergenca algoritma pri linearno ločljivih podatkih (levo) in linearno neločljivih podatkih (desno)



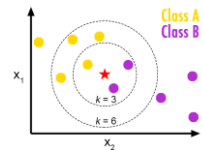
Pregled



- strojno učenje
 - uvod v strojno učenje
 - učenje odločitvenih dreves
 - učenje dreves iz šumnih podatkov (rezanje dreves)
 - regresijska drevesa
 - linearni modeli
 - metoda k najbližjih sosedov
 - ocenjevanje učenja

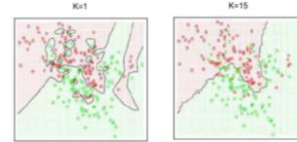
Metoda k najbližjih sosedov

- angl. k nearest neighbors
- lastnosti:
 - neparametrična** metoda (ne ocenjuje parametrov izbranega modela)
 - učenje na podlagi **posameznih primerov** (angl. *instance-based learning*)
 - leno učenje** (angl. *lazy learning*): z učenjem odlašajo vse do povpraševanja o novem primeru
- ideja: ob vprašanju po vrednosti odvisne spremenljivke za novi primer:
 - poišči k **primerov**, ki so **najbližji** glede na podano **mero razdalje**
 - napoved
 - pri klasifikaciji: npr. večinski razred med sosedi
 - pri regresiji: npr. povprečno vrednost/mediano označb sosedov
- v izogib neodločenemu glasovanju za večinski razred pri klasifikaciji običajno izberemo, da je k liho število



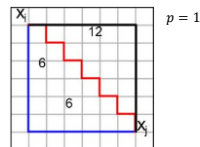
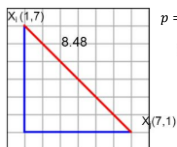
Metoda k najbližjih sosedov

- pomembna je izbira ustreznega k :
 - premajhen k : pretirano prilagajanje
 - prevelik k : prešibko posploševanje (pri $k = N$: napoved večinskega razreda)
 - v praksi običajno: $k = 5$



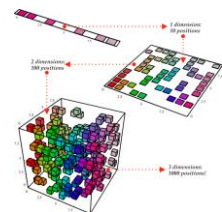
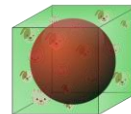
Metoda k najbližjih sosedov

- razdaljo običajno merimo z razdaljo Minkowskega: $L^p(x_i, x_j) = \left(\sum_k |x_{i,k} - x_{j,k}|^p \right)^{\frac{1}{p}}$
 - za $p = 2$ je to evklidska razdalja: $L^2(x_i, x_j) = \sqrt{\sum_k (x_{i,k} - x_{j,k})^2}$
 - za $p = 1$ je to manhattanska razdalja: $L^1(x_i, x_j) = \sum_k |x_{i,k} - x_{j,k}|$
- za zvezne atribute: razlika med vrednostima atributov (normalizacija?)
- za diskretne atribute: Hammingova razdalja (število diskretnih diskretnih atributov z ujemajočimi vrednostmi pri obeh primerih)



Opombe

- vpliv intervala vrednosti na izračunano razdaljo vpliva na najdene najbližje sosedse → potrebna **normalizacija**
- pri velikem številu dimenzij lahko postanejo primeri zelo oddaljeni – **prekletstvo dimenzionalnosti** (angl. *the curse of dimensionality*)
- implementacije iskanja najbližjih sosedov: $O(N)$, $O(\log N)$, $O(1)$



Pregled



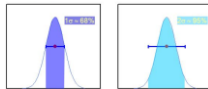
- strojno učenje
 - uvod v strojno učenje
 - učenje odločitvenih dreves
 - učenje dreves iz šumnih podatkov (rezanje dreves)
 - regresijska drevesa
 - linearni modeli
 - metoda k najbližjih sosedov
 - ocenjevanje učenja

Ocenjevanje učenja



- kriteriji za ocenjevanje hipotez:
 - točnost (angl. *accuracy*)
 - kompleksnost (angl. *complexity*)
 - razumljivost (angl. *comprehensibility*) – subjektivni kriterij
- ocenjevanje točnosti:
 - na **učnih** podatkih (angl. *training set, learning set*)
 - na **testnih** podatkih (angl. *testing set, test set*)
 - izločimo del učnih podatkov, s katerimi simuliramo ne-videne podatke
 - želim si, da je testna množica reprezentativna za nove podatke
 - uporabimo lahko **intervale zaupanja** v oceno uspešnosti na testni množici, ki upoštevajo število testnih primerov
 - na **novih** (ne-videnih) podatkih (angl. *new data, unseen data*)
 - na njih bo naučeni sistem dejansko deloval

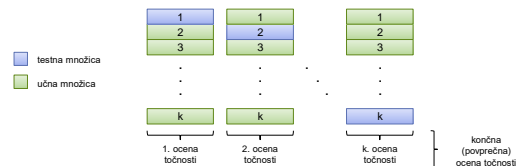
Ocenjevanje učenja



- nasprotujoča si cilja:
 - potrebujemo čim več podatkov za **uspešno učenje**
 - potrebujemo čim več podatkov za **zanesljivo ocenjevanje točnosti** (večje število testnih primerov nam daje ožji interval zaupanja v oceno točnosti)
- rešitev:
 - kadar je učnih podatkov dovolj, lahko izločimo **testno množico** (angl. *holdout test set*)
 - alternativa: **večkratne delitve** na učno in testno množico
- različni načini **vzorčenja testnih primerov**:
 - naključno, nenaključno (npr. prečno preverjanje)
 - poljubno ali stratificirano (zagotovimo enako porazdelitev razredov kot v učni množici)

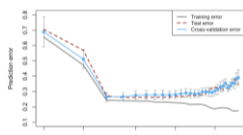
Prečno preverjanje

- poseben primer večkratnega učenja in testiranja
- k-kratno prečno preverjanje (angl. *k-fold cross-validation*):
 - celo učno množico razbij na k disjunktnih podmnožic
 - za vsako od k podmnožic:
 - uporabi množico kot testno množico
 - uporabi preostalih $k-1$ množic kot učno množico
 - povpreči dobjenih k ocen točnosti v končno oceno



Prečno preverjanje

- v praksi najpogostejše: $k=10$ (10-kratno prečno preverjanje)
- vplive izbranega razbitja podatkov na podmnožice lahko zmanjšamo tako, da tudi prečno preverjanje večkrat (npr. 10x) ponovimo (torej $10 \times 10 = 100$ izvajanj učnega algoritma) in rezultate povprečimo
- poseben primer prečnega preverjanja je metoda **izloči enega** (angl. *leave-one-out, LOO*)
 - k je enak številu primerov (vsaka testna množica ima samo en primer)
 - najbolj stabilna ocena glede učinkov razbitja na podmnožice
 - časovno zelo zamudno, primerno za manjše množice
- iz meritev na vseh podmnožicah je možno izračunati tudi varianco/ intervale zaupanja



Naivni Bayesov klasifikator