

1. ODKRIVANJE SKUPIN

16/1/20

Euklidova razdalja : $d(a,b) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$

Manhattanova razdalja : $d(a,b) = \sum_{i=1}^n |a_i - b_i|$

Minkowska razdalja : $d(a,b) = \left(\sum_{i=1}^n |a_i - b_i|^p \right)^{\frac{1}{p}}$

Razdalje med skupinami :

Single linkage

Complete linkage

Average linkage

Wardova razdalja (centroidi) :

$$d_w(C_u, C_v) = \sum_{x \in C_{uv}} d(x, R_{uv})^2 = \left(\sum_{x \in C_u} d(x, R_u)^2 + \sum_{x \in C_v} d(x, R_v)^2 \right)$$

2. METODA VODITELJEV

Izbira k groč \rightarrow minimizacija : $\sum_{i=1}^k \sum_{x \in C_i} d(v^{(i)}, x)$

preved
 \hookrightarrow SSE (sum of squared errors) : $SSE = \sum_{i=1}^k \sum_{x \in C_i} (v^{(i)} - x)^2$

Kohenzija: ~~oddaljenost~~

\hookrightarrow podobnost primerov znotraj skupine

Ločljivost:

\hookrightarrow oddaljenost primerov med različnimi skupinami

Silhuetni koeficient == silhueta razbitja (mera razbitja združuje)

silhueta - mera ocenjevanja
kvalitete razvrstitve

$$S = \frac{b-a}{\max(b,a)} \quad [-1,1] \rightarrow 0 \text{ je dobra silhueta točka}$$

Normalizacija :
$$z_{ij} = \frac{x_{ij} - \mu_j}{\sigma_j}$$

$$\sigma_j = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_{ij} - \mu_j)^2}$$

3. RAZVRŽEVANJE BESEDIL

16/1/20

inverse document frequency:

$$IDF(t) = \log \frac{|D|}{|\{d: t \in d\}| + 1}$$

→ št. dokumentov

→ št. el., ki vsebujejo el. t (kjer je $tf(t, d) \neq 0$)

$$tf-idf(t, d) = tf(t, d) \times idf(t)$$

Kosinusna podobnost:

$$a \cdot b = \|a\| \cdot \|b\| \cdot \cos \theta$$

$$\text{sim}(x, y) = \cos(\theta) = \frac{x \cdot y}{\|x\| \cdot \|y\|}$$

Podobnost po jaccardu: $J(x, y) = \frac{|x \cap y|}{|x \cup y|}$

4. PROJEKCIJE IN ZMANJŠEVANJE DIMENZIONALNOSTI PODATKOV

središče projekcije na enotski vektor:

$$\bar{x} = \frac{1}{n} \sum_i x^{(i)}$$

Razpršenost podatkov na projekciji:

$$\text{Var}(u_1^T X^T) = \frac{1}{n} \sum_{i=1}^n (u_1^T x^{(i)} - u_1^T \bar{x})^2$$

→ želimo maksimirati, zbrat primerni enotski vektor \vec{u}_1
 $\boxed{u_1^T \cdot u_1 = 1}$

(Numerična) potlačena metoda

→ računa se prvih n lastnih vrednosti pri velikem št. atributov
 → 1 naključno

Večrazsežnostno loščevanje: vizualizira podatke v obliki zemljevida (MDS) → pomanjkljivost: optimizira (tudi) oddaljene primere

Stohastična vložitev sosedov (SNE)

→ izboljšava MDS, poudarek je na sosedih, skoncentrira se na lokalnost
 → naključnem začetek vložitve

Kovariacijska matrika:

$$S = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^T$$

→ n. atributov → $\mathbb{R}^{n \times n}$

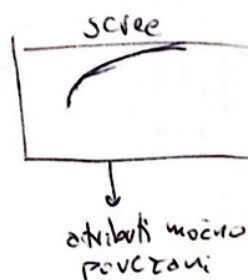
$$\text{Var}(u_1^T X^T) = u_1^T S u_1$$

$$S = (x_i - \bar{x}) \cdot (x_i - \bar{x})^T$$

$$\boxed{} = \boxed{} \cdot \boxed{}$$

$$S u_1 = \lambda_1 \cdot u_1$$

$$\text{Var}(u_1^T x^T) = \lambda_1$$



5. LINEARNA REGRESIJA

Kriterijska funkcija: $J = \frac{1}{2m} \sum_{i=1}^m (f(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m \epsilon^{(i)2}$

↳ jo minimiziramo

$$y = f(x) = ax + b$$

$$\frac{\partial J(a,b)}{\partial a} = \frac{1}{m} \sum_{i=1}^m ((ax^{(i)} + b) - y^{(i)}) x^{(i)}$$

$$\frac{\partial J(a,b)}{\partial b} = \frac{1}{m} \sum_{i=1}^m ((ax^{(i)} + b) - y^{(i)})$$

Multivariantna lin. reg.

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$$(y = ax + b) \\ = \theta_1 x + \theta_0$$

$$h(x) = \theta^T x$$

↳ prvi element je 1

Popravek za 1 parameter:

$$\theta_i \leftarrow \theta_i - \frac{\alpha}{m} (h_{\theta}(x) - y) x_i$$

za vse primere pa:

$$\theta \leftarrow \theta - \frac{\alpha}{m} \sum_{j=1}^m (h_{\theta}(x^{(j)}) - y^{(j)}) x_i^{(j)}$$

Normalizacija: \leftarrow avg

$$x_i \leftarrow \frac{x_i - \mu_i}{\sigma_i}$$

σ_i \leftarrow variance

$$\min_{\theta} \sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

kriterijska fun. =

$$J(\theta) = \frac{1}{2m} \sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

↳ za poln da ni odvisno od št. podatkov

Paketni pristop (za vse primere naenkrat):

$$\theta \leftarrow \theta - \alpha X^T (X\theta - \mathbf{y})$$

stohastični pristop (v for loopu vrak primer naenkrat)

6. REGULARIZACIJA

→ minimizira vsoto kvadratov teta

16/1/20

→ pomaga pri preprečevanju prevelikega prilaganja, nočem prevelikih koeficientov

$$RMSE(\mathcal{Y}) = \sqrt{\frac{\sum_{i=1}^k y^{(i)} - \hat{y}^{(i)}}{k}}$$

testna množica

→ napaka neodvisna od
prej domena odvisne sprem.

$$R^2(\mathcal{Y}) = 1 - \frac{\sum_k (y^{(i)} - \hat{y}^{(i)})^2}{\sum_k (y^{(i)} - \bar{y})^2}$$

delež razločene variance

avg

dobro model ima to 0, torej
je $R^2 \approx 1 \rightarrow \text{good!}$

$$J = \frac{1}{2m} \sum_{i=1}^m (f(x^{(i)}) - y^{(i)})^2 + \eta \sum_{j=1}^n \theta_j^2$$

cenovna fun.
za lin. reg.

želimo čim manjši
(parametre)

η → stopnja regularizacije
(0,1 ; 0,01 ali 0,001)

Gradient (odvod cenovne funkcije):

$$\frac{\partial}{\partial \theta_i} J(\theta) = \frac{1}{m} (h_{\theta}(x) - y) x_i + 2\eta \theta_i$$

L naredi binaren
nabor znakov
(jih nastavi na 0)

L2 nastavi neparametre
tete (parametre)
na zelo malo in ne 0.

$$\theta_j \leftarrow \theta_j \left(1 - \frac{\lambda}{m}\right) - \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_{ij}$$

7. LOGISTIČNA REGRESIJA

Log fun: $g(z) = \frac{1}{1 + e^{-z}}$

$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$

odvod: $\frac{dg(z)}{dz} = g(z) \cdot [1 - g(z)] = \frac{1}{1 + e^{-z}} \cdot \left(1 - \frac{1}{1 + e^{-z}}\right)$

$$p(y|x; \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}$$

Verjetje L

Verjetnost $L(\theta) = p(\bar{y}|X; \theta)$

$$= \prod_{i=1}^m h_{\theta}(x^{(i)})^{y_i} \cdot (1 - h_{\theta}(x^{(i)}))^{1-y_i}$$

iščemo parametre θ , da bo L največja

logaritem verjetja $\ell(\theta)$, ker je odvajanje
funkcij & mnogo produkti, precej enostavneje
logaritem verjetja (log likelihood):

$$\ell(\theta) = \log L(\theta)$$

$$= \sum_{i=1}^m [y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)}))]$$

gradient log. verjetja → maksimizira

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = \sum_{i=1}^m (y_i - h_{\theta}(x^{(i)})) x_{ij}$$

optimizacija parametrov modela: $\theta_j \leftarrow \theta_j + \alpha \sum_{i=1}^m (y_i - h_{\theta}(x^{(i)})) x_{ij}$

8. KLASIFIKACIJSKA DREVESA & GOZDOVI 17/1/20

Residualna entropija

$$H_{res} = H(C|X) = \sum_i p(x_i) H(C|x_i)$$

Information gain ratio:

$$IGR(x) = \frac{IG(x)}{H(x)} \rightarrow \text{uravnotežen informacijski prispevek}$$

Information gain:

$$IG(x) = H(C) - H(C|x)$$

↓
rez. entri.

koliko informacije nam prispeva poznavanje vrednosti posameznega atributa.

Mera nečistoče po Giniju:

$$Gini(C) = \sum_i p(C=c_i) \times (1 - p(C=c_i)) = 1 - \sum_i [p(C=c_i)]^2$$

pove: kako pogosto bi bil za naključno izbran primer napakno napovedan razred, ki bi ga uteženo naključno napovedali

nič je, ko vsi ~~primeri~~ primeri pripadajo istemu razredu

Kategorizacija / diskretizacija atributa

Gozdovi:

Raw importance: $RI(x) = \frac{1}{k} \sum_{i=1}^k C_i - C_{ix}$
(pomembnost atributa)

9. PRIPOROČILNI SISTEMI

17/1/20

napoved preferenčne ocene uporabnika u za izdelek i : $\hat{r}_{ui} = \frac{\sum s(u, u') \times r_{u'i}}{\sum s(u, u')}$

↑ uteži podobnosti med uporabnikoma

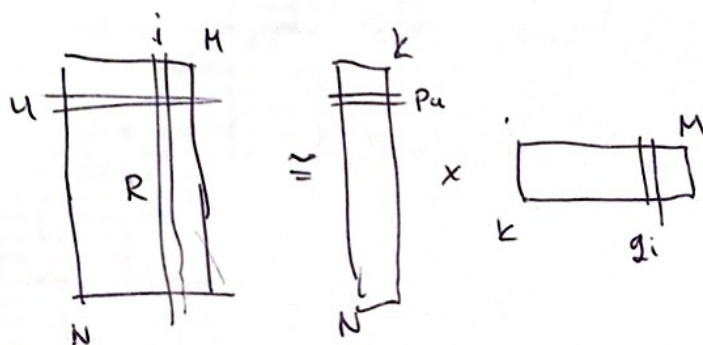
↓ vsota uteži (normalizacija)

podobnost: (cos)

$$S_c(u, u') = \frac{r_u \cdot r_{u'}}{|r_u| |r_{u'}|}$$

$P_{N \times K}$ in $Q_{K \times M} \rightarrow R \approx PQ$
 $k \ll M, N$
 $k \dots$ stopnje razcepa

Latentni profil uporabnika p_u
 l.p. stvari: g_i



približek ocene:
 $\hat{r}_{ui} = P_u^T \cdot g_i$

$$e_{ui}^2 = \frac{1}{2} \left(r_{ui} - \sum_{k=1}^k P_{uk} g_{ki} \right)^2$$

odvod $\frac{\partial e_{ui}^2}{\partial P_{uk}} = -e_{ui} \cdot g_{ki}$

$\frac{\partial e_{ui}^2}{\partial g_{ki}} = -e_{ui} \cdot P_{uk}$

update:

$P_{uk} \leftarrow P_{uk} + \alpha e_{ui} g_{ki}$

$g_{ki} \leftarrow g_{ki} + \alpha e_{ui} P_{uk}$

BMF:

inkrementalna simultana matricna faktorizacija
 (iskanje razcepa matrice na P in Q)

vektorsko:

$P_u \leftarrow P_u + \alpha e_{ui} g_i$

$g_i \leftarrow g_i + \alpha e_{ui} P_u$

10. POVEZOVALNA PRAVILA

17/1/20

$\sigma(X)$... št. transakcij, ki jih vsebuje X

$s(X)$... delež podprtih transakcij aka
delež transakcij, ki vsebujejo X
↑
vseh

$$s(X) = \frac{\sigma(X)}{|T|} = \frac{\sigma(X)}{N}$$

Teorem 1

$$s(X) \geq \text{minsupp} \Rightarrow s(Y) \geq \text{minsupp}, Y \subseteq X$$

Teorem 2

$$s(X) < \text{minsupp} \Rightarrow s(Y) < \text{minsupp}, X \subseteq Y$$

Teorem 3: podpora uabora nikoli ne presega njegove podmnožice
↑
podpore

$$\forall X, Y \in \mathcal{L} : X \subseteq Y \Rightarrow s(Y) \leq s(X)$$

Apriori algoritem

Povezovalna pravila

Podpora (delež transakcij, kjer
uže poznamo najdeno vse stvari
iz povezovalnega pravila) : $s(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{N}$ minsupp

Zaupanje (delež transakcij, ki
vsebujejo desno stran
pravila Y med transakcijami,
ki vsebujejo levo stran pravila X) : $c(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}$ minconf