

Poslovna inteligenca

1. izpitni rok

25. januar 2016

Priimek in ime (tiskano): _____

Vpisna številka: _____

Naloga	1	2	3	4	5	Vsota
Vrednost	7	3	5	6	6	27
Točk						

1. Kriterijska funkcija, ki jo želimo maksimizirati pri logistični regresiji, je

$$l(\Theta) = \sum_{i=1}^m y^{(i)} \log h_{\Theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\Theta}(x^{(i)}))$$

kjer je funkcija $h_{\Theta}(x) = g(\Theta^T x)$ logistična funkcija linearne kombinacije vhodnih spremenljivk (atributov). Z uporabo metode gradientnega spusta lahko izpeljemo pravilo za iterativni popravek i -tega parametra linearne kombinacije:

$$\Theta_j \leftarrow \Theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\Theta}(x^{(i)})) x_j^{(i)}$$

Problem opisanega postopka je preveliko prileganje uĉnim podatkom. Zato uvedemo regularizacijo.

- [1] (a) Kako vpliva regularizacija na vrednost parametrov Θ ?
- [1] (b) Ali je toĉna trditev: veĉja je stopnja regularizacije, manjĹa je klasifikacijska toĉnost na uĉnih podatkih?
- [1] (c) Ali je toĉna trditev: veĉja je stopnja regularizacije, veĉja je klasifikacijska toĉnost na testnih podatkih?
- [2] (d) V zgornjo enaĉbo za kriterijsko funkcijo $l(\Theta)$ dodaj ĉlen z regularizacijo (uporabi tako enaĉbo oziroma tako regularizacijo, ki jo boĹ znal odvajati).
- [2] (e) Kako se z regularizacijo spremeni enaĉba za iterativni popravek? ZapiĹi novo enaĉbo popravka, ki upoĹteva regularizacijo. (Ne priĉakujemo, da znaĹ enaĉbo na pamet. Œe najbolj enostavno boĹ reĹitev dobil z odvodom kriterijske funkcije).

Pri odgovorih skuĹaj upoĹtevat, da je med parametri Θ parameter Θ_0 uporabljen kot konstantni ĉlen v linearni funkciji h_{Θ} .

Solution:

- Regularizacija zmanjĹa vrednosti parametrov, predvsem tistih, ki bi bili brez regularizacije visoki.
- Da.
- Ne.
- Ker $l(\Theta)$ maksimiziramo, Źelimo pa, da bi bili parametri ĉim manjĹi, dodamo ĉlen

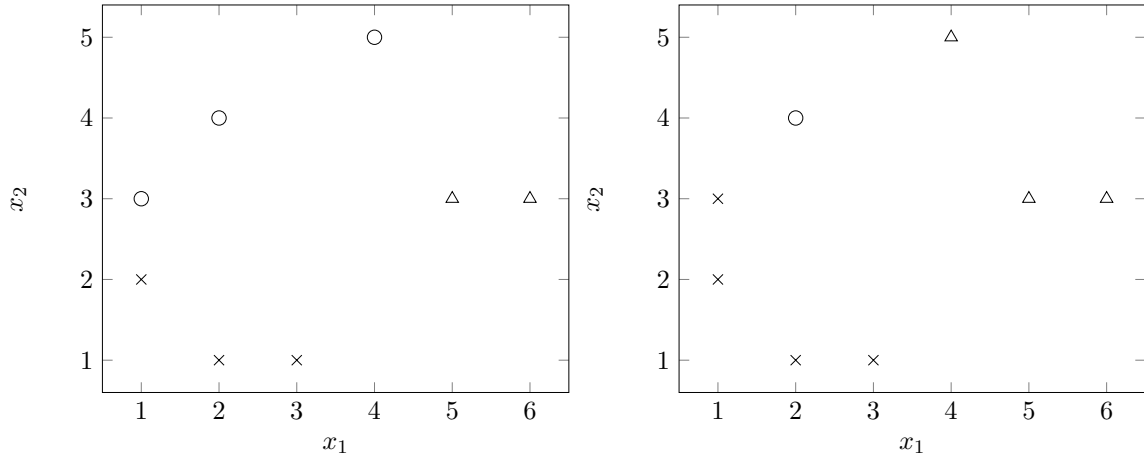
$$-\frac{\lambda}{2} \sum_{j=1}^n \Theta_j^2$$

- Namesto Θ_j na desni strani imamo $\Theta_j(1 - \alpha\lambda)$

- [3] 2. Časopisno hišo, ki objavlja novice na spletnih straneh, zanima model, ki bi na podlagi besedila novice ocenil, ali bo ta dobro brana. Za naš pilotni projekt so nam pripravili zbirko 10.000 novic in pri vsaki označili, ali je bila dobro ali slabo brana. Odločili smo se, da bomo za potrebe modeliranja novice predstavili z vektorjem prisotnosti besed. Vseh 10.000 novic skupaj uporablja 13.345 različnih besed. Da bi zadevo poenostavili, smo zato izbrali manjši nabor 1.000 besed tako, da smo vsako predstavili kot atribut (prisotnost besede v novici), stopnjo povezanosti atributa z razredom pa ocenili na podlagi informacijskega prispevka. Izbrali smo 1.000 besed z najvišjim informacijskim prispevkom. Na tako dobljeni podatkovni množici (10.000 novic, vsaka opisana z vektorjem prisotnosti 1.000 besed) smo potem ovrednotili uporabo logistične regresije ter na prečnem preverjanju izmerili AUC, ki je znašal 0.95. Časopisno hišo smo obvestili, da smo na njihovem vzorcu dobili izjemno visoko točnost in da je logistična regresija primerna metoda za gradnjo modelov za napovedovanje branosti novic.

Komentiraj primernost izbora postopkov ter upravičenost našega zaključka. Če se s kakšnim delom opisanega postopka ne strinjaš, predlagaj alternativno rešitev.

- [5] 3. Dobili smo podatke, ki so opisani z dvema numeričnima atributoma (x_1 in x_2) in jih zato lahko enostavno predstavimo v kartezični ravnini. Na podatkih smo uporabili dve različni tehniki razvrščanja v skupine. Rezultate razvrščanja za prvo metodo (metoda A) prikazuje prva slika, za drugo (metoda B) pa druga slika (pripadnost skupini je označena s simbolom, s katerim smo izrisali točke, ki ponazarjajo primere. Prva metoda je tako identificirala dve skupini s po tremi primeri - križci in krožci - in eno skupino z dvema primeroma - trikotnika).



Rezultate razvrščanja metod A in B oceni s tehniko silhuete, in sicer tako, da za vsako od obeh razvrstitev izračunaš pripadajočo (približno) vrednost silhuete. Pri računanju silhuete uporabljaš Evklidske razdalje. Katera razvrstitev je boljša in če, je razlika med razvrstitvama velika? (Namig: uporabi ravnilo, računanje oziroma oceno parametrov čimbolj poenostavi, točnost računanja pa omeji na največ eno decimalno mesto. Rezultate podajaj jasno, s tabelo.)

$$s_i = \frac{b_i - a_i}{\max a_i, b_i}$$

Solution:

x1	x2	a	b	s	a	b	s
3.0	1.0	1.5	3.0	0.5	2.0	3.2	0.4
2.0	1.0	1.8	3.0	0.4	1.5	3.0	0.5
1.0	2.0	1.3	2.5	0.5	1.5	2.2	0.3
1.0	3.0	2.5	2.0	-0.3	2.0	1.4	-0.4
2.0	4.0	1.8	2.7	0.4	0.0	1.4	1.0
4.0	5.0	2.8	2.4	-0.1	2.5	2.2	-0.1
5.0	3.0	1.0	2.8	0.6	1.5	3.0	0.5
6.0	3.0	1.0	3.8	0.7	1.7	3.8	0.5
average				0.3			0.3

obe razvrstitvi imata zelo podobno silhueto

Stran je prazna, da lahko nanjo rešujete nalogo.

4. Pet prijateljev je v spodnji tabeli dejavnosti označilo z ocenami od 1 do 5.

	tango	salsa	plezanje	bordanje	plavanje	kolesarjenje
Agnieszka			5		4	3
Marko	5		4	4	3	3
Nejc		1	2		3	2
Sara	2	1		5		5
Urška		2	4	2		1

Pomagajte Agnieszki, da se odloči med tangom, salso in bordanjem.

- [3] (a) Agnieszkiine potencialne aktivnosti razvrstite z metodo, ki deluje na podlagi podobnosti med prijatelji.
- [3] (b) Agnieszkiine potencialne aktivnosti razvrstite z metodo, ki deluje na podlagi podobnosti med aktivnostmi.

$$s_c(u, u') = \frac{\mathbf{r}_u \mathbf{r}_{u'}}{|\mathbf{r}_u| |\mathbf{r}_{u'}|}$$

$$s_J = \frac{\|\mathbf{r}_u \cap \mathbf{r}_{u'}\|}{\|\mathbf{r}_u \cup \mathbf{r}_{u'}\|}$$

Solution:

(-1 točka če ne množim podobnosti z oceno)

(-0.5, če je normaliziral tako, da ni seštev vseh razdalj - tudi tistih brez ocene)

a) Kosinusne razdalje:

Če predpostavimo neznane razdalje = 0 -> sicer slaba predpostavka (-0.5 točke pri delu a)

AM 0.67

AN 0.93

AS 0.29

AU 0.65

Ni treba normalizirati, ker je isto!

tango = 3.93

salsa = 2.52

bord = 5.43

POSEBNOSTI

- če je kdo opazil, da je Marko Agnieszki najbliže in je dal Markove ocene, je dobil 2

b) Pri B potrebujem normalizacijo. Če manjka normalizacija je (-1 točka).

Stran je prazna, da lahko nanjo rešujete nalogo.

5. Podana je tabela dobičkov, ki zajema tri stanja (S_i) in tri alternative (a_j). Tabela vključuje verjetnosti nastopa posameznih stanj.

Stanje	Verjetnost $p(S_i)$	Alternative		
		a_1	a_2	a_3
S_1	0,2	150	180	130
S_2	0,5	190	160	140
S_3	0,3	120	150	170

- [2] (a) Izračunajte pričakovane koristnosti za vse alternative. Za katero alternativo bi se odločili?
- [1] (b) Kako bi se odločili po kriteriju optimista, če verjetnosti nastopa posameznih stanj ne bi poznali?
- [1] (c) Kako bi se odločili po kriteriju pesimista, če verjetnosti nastopa posameznih stanj ne bi poznali?
- [2] (d) Kako bi se odločili po Hurwitzovem kriteriju, če verjetnosti nastopa posameznih stanj ne bi poznali in bi za vrednost koeficienta tveganja d vzeli 0,3?

Solution:

$$a) u(a_1) = 0,2 * 150 + 0,5 * 190 + 0,3 * 120 = 161$$

$$u(a_2) = 0,2 * 180 + 0,5 * 160 + 0,3 * 150 = 161$$

$$u(a_3) = 0,2 * 130 + 0,5 * 140 + 0,3 * 170 = 147$$

Odločili bi se za varianto a_1 ali a_2 , ki imata enakovredno najugodnejšo pričakovano koristnost.

Če omeni samo eno, dam 1.5 točke.

b) Odločili bi se za varianto a_1 .

c) Odločili bi se za varianto a_2 .

$$d) u(a_1) = d * \max(150;190;120) + (1-d) * \min(150;190;120) = 0,3 * 190 + 0,7 * 120 = 144$$

$$u(a_2) = d * \max(180;160;150) + (1-d) * \min(180;160;150) = 0,3 * 180 + 0,7 * 150 = 159$$

$$u(a_3) = d * \max(130;140;170) + (1-d) * \min(130;140;170) = 0,3 * 170 + 0,7 * 130 = 142$$

Odločili bi se za varianto a_2 .

Če je zamešal kam paše d je 169 171 158 -> a_2 . Tudi priznam.