

Predicting League of Legends match outcome with Bayesian methods

Z. Nagelj

April 21, 2019

1 Uvod

Logisticna regresijska se uporablja za analizo povezavo med kategoricno odvisno spremenljivko in neodvisnimi spremenljivkami. Poleg logisticne regresije se za analizo kategoricnih spremenljivk uporablja diskriminantna analiza, ki za razliko od logisticne regresije predpostavlja normalno porazdelitev neodvisnih spremenljivk.

2 Logisticna regresija

Kot vhodni podatek za logisticno regresijo dobimo podatkovni set N točk. Vsaka i -ta točka sestavlja set m -tih neodvisnih spremenljivk in kategoricna odvisna spremenljivka Y_i z dvema možnima izidoma.

2.1 Logit in logisticna transformacija

Naši kategoricni odvisni spremenljivki najprej dodelimo numerične vrednosti (0 in 1). Povprečje na vzorcu predstavlja delež ugodnih izidov p , razmerje $p/(1-p)$ pa obeti (odds). Logit transformacijo definiramo kot logaritem obetov (log odds):

$$l = \text{logit}(p) = \log \frac{p}{1-p}$$

S pomočjo te transformacije preidemo iz omejene zaloge vrednosti p na intervalu $[0, 1]$, na obete $p/(1-p)$ omejene z zalogo vrednosti $[0, \infty)$ in na koncu na logaritem obetov z zalogo vrednosti $(-\infty, \infty)$. Inverzno transformacijo imenujemo logisticna transformacija:

$$p = \text{logistic}(l) = \frac{\exp l}{1 + \exp l}$$

S transformacijo se izognemo problemu omejene zaloge vrednosti odvisne spremenljivke. Potencialno bi lahko izbrali tudi kakšno drugo transformacijo (probit).

2.2 Logisticni model

Kategoricno odvisno spremenljivko definiramo kot slučajno spremenljivko Y_i porazdeljeno po Bernoulliju s pričakovano vrednostjo p_i . Vsak izid je torej določen s svojo neznano verjetnostjo p_i , ki je določena na podlagi neodvisnih spremenljivk.

$$Y_i|x_{1,i}, \dots, x_{m,i} \sim \text{Bernoulli}(p_i)$$

$$E[Y_i|x_{1,i}, \dots, x_{m,i}] = p_i$$

$$P(Y_i = y|x_{1,i}, \dots, x_{m,i}) = p_i^y(1 - p_i)^{(1-y)}$$

Ideja je zelo podobna kot pri linearni regresiji, torej verjetnost p_i modeliramo kot linearno kombinacijo neodvisnih spremenljivk. Razlika je v tem, da verjetnosti transformiramo s pomočjo logit funkcije. V modelu nastopa dodaten intercept člen, zato imamo $m + 1$ regresijskih koficientov β .

$$\text{logit}(E[Y_i|X_i]) = \text{logit}(p_i) = \log \frac{p_i}{1 - p_i} = X_i\beta$$

Oziroma:

$$E[Y_i|X_i] = p_i = \text{logit}^{-1}(X_i\beta) = \frac{1}{1 + \exp^{-X_i\beta}}$$

$$P[Y_i = y_i|X_i] = p_i^{y_i}(1 - p_i)^{1-y_i} = \frac{\exp^{y_i X_i \beta}}{1 + \exp^{X_i \beta}}$$

2.3 Dolocanje vrednosti regresijski koficientov

Regresijski koficienti in verjetnosti p_i so dolocene optimizacijo, na primer MLE. Za lažjo predstavbo si najprej ogledamo MLE za enostaven primer Bernoullija:

2.4 MLE za Bernoulli(p)

Zapišemo enačbo za verjetje in jo logaritmiramo.

$$L = \prod_{i=1}^n p^{Y_i}(1 - p)^{1-Y_i} = p^{\sum_{i=1}^n Y_i} (1 - p)^{n - \sum_{i=1}^n Y_i}$$

$$l = \log(L) = \sum_{i=1}^n Y_i \log(p) + (n - \sum_{i=1}^n Y_i) \log(1 - p)$$

Cenilko za \hat{p} dolocimo s prvim parcialnim odvodom, ki ga enacimo z 0. Za asimptotski interval zaupanja cenilke dolocimo Fisherjevo informacijsko matriko. Saj populacijske vrednosti p ne poznamo Fisherjevo informacijo dolocimo z oceno \hat{p} .

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$I(p) = E\left[-\frac{\partial^2 l}{\partial p^2}\right] = \frac{1}{p(1-p)}$$

$$(\hat{var}) = \frac{1}{nI(\hat{p})} = \frac{\hat{p}(1-\hat{p})}{n}$$

$$(\hat{SE}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Za velik n je naša cenilka porazdeljena približno normalno:

$$\hat{p} \sim_{CLI} \text{Normal}\left(p, \frac{p(1-p)}{n}\right)$$

2.5 MLE za Bernoulli(p_i)

Saj ima pri logisticne regresije vsak izid svojo verjetnost p_i je naša enacba za verjetje naslednja:

$$L = \prod_{i=1}^n p_i^{Y_i} (1 - p_i)^{1-Y_i} = p_i^{\sum_{i=1}^n Y_i} (1 - p_i)^{n - \sum_{i=1}^n Y_i} = \left(\frac{\exp^{X_i \beta}}{1 + \exp^{X_i \beta}} \right)^{\sum_{i=1}^n Y_i} \left(\frac{1}{1 + \exp^{X_i \beta}} \right)^{n - \sum_{i=1}^n Y_i}$$

$$l = \log(L) = \sum_{i=1}^n Y_i \log \left(\frac{\exp^{X_i \beta}}{1 + \exp^{X_i \beta}} \right) + \left(n - \sum_{i=1}^n Y_i \right) \log \left(\frac{1}{1 + \exp^{X_i \beta}} \right)$$

Posamezne regresijske koficiente dobimo s parcialnim odvodom logaritma verjetja. Saj rešitev v zaključeni formi ne obstaja uporabimo numerice metode (npr. Newtonova metoda). Prav tako določimo informacijsko matriko za določanje asimptotske kovariančne matrike in intervalov zaupanja.

3 Statisticno testiranje regresijskih koficientov

Za testiranje hipotez nenicelnosti regresijski koficientov sta v uporabi Waldov test in test razmerja verjetij.

3.1 Waldov test

Waldov test se uporablja za določanje statistične značilnosti posameznih regresijskih koficientov (podobno kot t-test pri linearni regresiji). Za koficiente pridobljene z MLE je testna statistika naslednja:

$$Z = \frac{\hat{\beta}_i}{\hat{SE}}$$

Nicelelna hipoteza testira $H_0 : \beta_i = 0$. Kot velja za vse cenilke pridobljene po metodi največjega verjetja so te asimptotsko nepristranske (dosledne) in normalno porazdeljene okoli prave vrednosti z varianco $\frac{1}{nI(p)}$.

3.2 Test razmerja verjetij

Pri testu nas zanima logaritem Wilksovega lambda pri dveh različnih modelih, pri čemer je en model gnezden (podmnožica drugega). Primerjali bomo polni model s \mathbf{k} regresijskimi koficienti in delni model z \mathbf{m} regresijskimi koficienti, kjer je $\mathbf{m} < \mathbf{k}$. Nicelna hipoteza (delni model) trdi da so testirani regresijski koficienti enaki 0. Alternativna hipoteza (polni model) pa trdi, da so vsi regresijski koficienti različni od 0. Pod nicelno hipotezeo torej testiramo nicelnost $\mathbf{k} - \mathbf{m}$ β . Naša testna statistika je:

$$LR = -2 \log \Lambda = -2 \log \frac{L(H_0)}{L(H_A)} = -2(l(H_0) - l(H_A)) = -2(l(\hat{\beta}^{(0)}) - l(\hat{\beta}))$$

Asimptotsko je ta porazdeljena s χ^2 z $\mathbf{k} - \mathbf{m}$ stopinjami prostosti. Test razmerja verjetij se od Waldovega testa razlikuje po tem, da je potrebno narediti dva modela pod različnimi hipotezami.

4 Naloga

4.1 Opis

Na vzorcu bolnikov z rakom primerjamo dve vrsti operacije. Zanima nas ali obstaja povezanost s stadijem bolnika in ali ena vrsta operacije povzroca manj zapletov. Cilj naloge je analizirati napako I. reda pri testiranju hipotez na nacin, da testiramo vsako spremenljivko posebj.

4.2 Postopek testiranje

Sledili bomo naslednjim korakom:

4.2.1 Stadij

1. Generiranje podatkov pod dano hipotezo
2. Izdelava 4 modelov logisticne regresije z informacijo o posameznem stadiju
3. Pridobimo p-vrednosti Waldovega testa vseh 4 modelov z delnimi podatki
4. Izdelava 1 modela logisticne regresije z informacijo o vseh stadijih (1 spremenljivka)
5. Pridobimo p-vrednosti testa razmerja verjetij za model z vsemi podatki
6. Analiziramo porazdelitve p-vrednosti in primerjamo deleže zavrženih hipotez

4.2.2 Zaplet

1. Generiranje podatkov pod dano hipotezo
2. Izdelava 10 modelov logisticne regresije z informacijo o posameznem zapletu
3. Izdelava 1 modela logisticne regresije z informacijo o vseh zapletih (10 spremenljivk)
4. Pridobimo p-vrednosti Waldovega testa vseh 11 modelov
5. Analiziramo porazdelitve p-vrednosti in primerjamo deleže zavrženih hipotez

4.3 Pricakovani rezultati

Zagotovo, pričakujemo, da bo naveden nacin testiranja pri stadiju napacen, saj so stadiji med seboj neodvisni, cesar v modelih ne zajamemo. Prav tako v model ne vkljucimo informacije o ostalih stadijih in jih obravnavamo kot enakovredne. Pravilno bi stadij tako, da ga v celoti vkljucimo v model ter z testom razmerja verjetij testiramo nenicelnost regresijskega koficienta.

Pri obravnavi zapletov operacij, ob predpostavki da so si te neodvisni, predvidevam, da tak nacin testiranja ne bi bil napacen. V realnost pa temu zagotovo ni tako in so posamezni zapleti med seboj povezani, zato bi zagotovo naleteli na enake probleme kot pri testiranju posameznega stadija.

5 Podatki

5.1 Generiranje

Saj je v nalogi določeno, da je vzorec velik 300 pacientov, kjer vsaka polovica prejme en tip operacije, najprej generiramo večji vzorec $k = 10000$ bolnikov iz katerega bomo vzorčili. Pridobiti moramo vrednosti spremenljivke stadij in desetih spremenljivk zaplet.

5.1.1 Stadij

Definiramo populacijske verjetnosti za vsak stadij (Tabela 1), ter glede na njih s funkcijo `sample` izžrebamo stadij vsakega bolnika in določimo modelsko matriko. Verjetnosti stadijev se morajo sešteti v 1, saj so med seboj odvisni.

	verjetnost
1	0.60
2	0.25
3	0.10
4	0.05

Table 1: Populacijski verjetnosti posameznega stadija raka

5.1.2 Zaplet

Definiramo populacijske verjetnosti za vsak zaplet (Tabela 2), ter glede na njih s funkcijo `rbinom` izžrebamo ali se je posamezen zaplet zgodil.

	verjetnost
Stadij 1	0.10
Stadij 2	0.12
Stadij 3	0.14
Stadij 4	0.16
Stadij 5	0.18
Stadij 6	0.20
Stadij 7	0.30
Stadij 8	0.40
Stadij 9	0.50
Stadij 10	0.60

Table 2: Populacijski verjetnosti posameznega zaplete pri operaciji

5.1.3 Vzorec

Ko imamo določene vrednosti spremenljiv na podlagi definirane hipoteze (oz. regresijskih koficientov) določimo linearne kombinacije spremenljivk ter s pomočjo logit transformacije verjetnost za tip operacije za vsakega izmed k bolnikov. Izbran tip operacij pridobimo iz Bernoullijeve porazdelitve, glede na p_i .

Da zadostimo specifikacijam naloge iz vzorca 10000 bolnikov ob vsaki izmed $m = 1000$ iteracij naključno izžrebamo 150 operacij vsakega tipa.

5.2 Hipoteze

Podatke generiramo glede na 2 hipoteze, ničelno in alternativno. Glede na ničelno hipotezo so vsi regresijski koficienti enaki 0. Vrednosti regresijski koficientov pod alternativno hipotezo so vidni v Tabeli 3.

	stadij	zaplet
beta0	0.00	0.00
beta1	0.00	-3.50
beta2	0.10	3.00
beta3	0.60	-2.50
beta4	-1.20	2.00
beta5		-1.50
beta6		1.00
beta7		-0.80
beta8		0.60
beta9		-0.40
beta10		0.20

Table 3: Vrednosti regresijski koficientov pri H_A glede na neodvisno spremenljivko