

Линеарни статистички модели
припрема 2019.

1. Посматрајмо модел линеарни модел $Y = X\beta + \varepsilon$ са слободном чланом различитим од нуле и дизајн матрицом X $n \times (p + 1)$, $\beta = (\beta_0, \dots, \beta_p)$ вектор непознатих параметара, и ε има вишедимензиону нормалну $\mathcal{N}_n(0, \sigma^2 I)$ расподелу, при чему је σ^2 такође непознат параметар.

а) Оценити вредност $C \times \beta$, где је $C = (c_0, \dots, c_p)$ методом максималне веродостојности. Коју расподелу има та оцена? Одредити 90% интервал поверења за $\beta_0 + \beta_1 + \dots + \beta_p$.

б) Скуп података `longley` (доступан у `R`-у) има 7 променљивих: `GNP.deflator`, `GNP`, `Unemployed`, `Armed.Forces`, `Population`, `Year` и `Employed`. Направити линеаран модел који описује зависност `Employed` од осталих променљивих. Одредити 95% интервал поверења за $\beta_{Population} - \beta_{GNP}$ и тестирати хипотезу $H_0 : \beta_{Population} = \beta_{GNP}$.

в) Шта представља мултиколинеарност? Да ли је потребно познавање зависне променљиве за детекцију исте? Одговор образложити. Како се тај проблем превазилази коришћењем назубљене регресије? Од чега зависи оцена непознатих параметара добијена овом методом? Коју расподелу има оцена?

г) Испитати да ли у моделу под б) постоји проблем са мултиколинеарношћу? Ако постоји, решити проблем избацивањем одговарајућих предиктора. Упоредити коефицијенте детерминације почетног и новог модела.

д) Показати да је $\sum_{i=1}^n e_i = 0$

ђ) Да ли су резидуали модела независни? А грешке модела? Одговор образложити.

е) На шта се односи услов хомоскедастичности у линеарном моделу? Како се детектује и како се може решити проблем хетероскедастичности?

ж) Направити линеарни модел који описује зависност променљиве `dist` од променљиве `speed` из скупа података `cars` (уграђен у `R`). У модел додати и квадратни члан за предиктор `speed`. Испитати графички да ли постоји проблем хетероскедастичности. Уколико постоји проблем, применити Бокс-Кокс трансформацију зависне променљиве. Да ли је проблем ублажен? Који модел (са или без трансформације) боље описује податке?

2. а) Објаснити зашто се регресиона функција зависне променљиве која има Бернулијеву расподелу не може моделирати линеарним моделом. б) Навести две линк функције које су погодне за моделирање Бернулијевих зависних променљивих.

б) Написати функцију веродостојности за логистички модел у коме је X дизајн матрица ред $n \times p + 1$, а вектор зависних променљивих Y .

в) Која је предност коришћења логит трансформације над онсталим линк функцијама?

г) Скуп података `leukemia` (доступан у пакету `LSMhelp`) садржи једну бинарну променљиву `REMISS` и непрекидне променљиве `CELL`, `SMEAR`, `INFIL`, `LI`, `BLAST` и `TEMP`. Направити модел логистичке регресије којим се описује утицај ових променљивих на променљиву `REMISS`.

д) Написати који сте добили модел уколико се посматрају сми могући предиктори.

ђ) Како изгледа Валдова тест статистика? На основу Валдовог теста, да ли су неки предиктори значајни и који?

е) Тестирати да ли постоји било какав утицај предиктора на зависну променљиву. Одговор образложити.

ж) Да ли постоји значајна разлика између модела који укључује само `LI` као предиктор и полазног модела?

з) Који од ових модела је бољи на основу AIC ?

е) Шта представља девијација модела? Да ли се у крајњем моделу који сте добили се девијација може искористити за тестирање да ли је модел адекватан?

```

1 ##### Zadatak 1.
2 ### 1.b)
3 model <- lm(Employed ~ ., longley)
4 summary(model)
5 # Pravimo interval poverenja za b_pop - b_gnp kao na petom casu
6 l <- c(0, 0, -1, 0, 0, 1, 0)
7 X <- model.matrix(model)
8 XtXi <- solve(t(X)%*%X)
9 n <- nrow(X)
10 p <- length(model$coefficients) - 1
11 estim <- t(1) %*% model$coefficients
12 stdev <- sqrt(sum(model$res^2)/(n-p-1) * t(1) %*% XtXi %*% 1)
13 lwr <- estim - stdev * qt(0.975, n-p-1)
14 upr <- estim + stdev * qt(0.975, n-p-1)
15 c(lwr, upr)
16
17 # Ovaj interval ukljucuje nulu, sto znaci da ne mozemo da odbacimo hipotezu da
18 # je b_pop = b_gnp
19 # Drugi nacin za testiranje H0: b_pop = b_gnp je koriscenjem linear_hypothesis
20 linear_hypothesis(model, t(c(0, 0, -1, 0, 0, 1, 0)), 0)
21 # ili anova...
22 anova(lm(Employed ~ . -GNP - Population + I(GNP + Population), longley), model)
23
24 ### 1.g) Multikolinearnost
25 # Mozemo da gledamo vif
26 library(car)
27 vif(model)
28 # Mnoge vrednosti su velike, dakle postoji prisustvo multikolinearnosti
29
30 # Promenljive Unemployed i Armed.Forces imaju vidno manji vif od ostalih, tako
31 # da cemo ih sacuvati.
32 # Pogledajmo korelacije ostalih promenljivih
33 cor(longley[, -c(3, 4, 7)])
34 # Primecujemo izuzetno velike korelacije, mozemo i nacrtati grafike
35 plot(longley[, -c(3, 4, 7)])
36 # Velika linearna zavisnost je prisutna, tako da mozemo ih sve zameniti samo
37 # jednom od promenljivih, npr. Year.
38 model2 <- lm(Employed ~ Armed.Forces + Unemployed + Year, longley)
39 summary(model2)
40 vif(model2)
41 # U ovom modelu su vif-ovi svi mali (<30) pa nemamo vise problem
42 # multikolinearnosti, dok je koeficijent determinacije smanjen sa 0.9955 na
43 # 0.9928, sto je vrlo mala razlika, pa je dobar smanjeni model.
44
45 ### 1.zh) Heteroskedasticnost
46 model <- lm(dist ~ speed + I(speed^2), cars)
47 # Ispitujemo prisustvo heteroskedasticnosti npr. grafikom zavisnosti apsolutnih
48 # vrednosti standardizovanih reziduala od ocenjenih vrednosti dist.
49 plot(abs(rstandard(model)) ~ fitted(model))
50 # Vidi se da disperzija raste sa porastom dist, tako da imamo problem
51 # heteroskedasticnosti.
52
53 # Trazimo odgovarajucu box-cox transformaciju
54 library(MASS)
55 bc <- boxcox(model)
56 lambda <- bc$x[which.max(bc$y)] # = 0.38
57 # Pravimo funkciju za transformaciju
58 BC <- function(y) {
59   (y^lambda - 1)/lambda
60 }
61 # Fitujemo model sa boxcox transformisanom zavisnom promenljivom
62 model2 <- lm(BC(dist) ~ speed + I(speed^2), cars)
63 summary(model2)
64 plot(abs(rstandard(model2)) ~ fitted(model2))
65 # Sada je heteroskedasticnost vidno ublazena.
66
67 ##### Zadatak 2.
68 library(LSMhelp) # instalacija: devtools::install_github("blaza/lsmhelp")
69 # ili preuzeti podatke: https://newonlinecourses.science.psu.edu/stat501/sites/onlinecourses.science.psu.edu.stat501/
70 #   files/data/leukemia_remission/index.txt
71 # Pravimo pocetni model sa svim prediktorima
72 m <- glm(REMISS ~ ., leukemia, family=binomial)
73 # U summary vidimo vrednosti Valdivog testa (zvezdice)
74 summary(m)
75 # Svi koeficijenti su beznacajni
76
77 # Testom kolicnika verodostojnosti proveravamo da li postoji uticaj prediktora
78 # na zavisnu. Poredimo model koji ukljucuje samo slobodan clan i nas model.
79 anova(glm(REMISS ~ 1, leukemia, family=binomial), m, test = "Chisq")
80 # p vrednost < 0.05, pa zakljucujemo da postoji znacajan uticaj
81
82 # Poredimo model samo sa LI sa nasim modelom
83 anova(glm(REMISS ~ LI, leukemia, family=binomial), m, test = "Chisq")
84 # p vrednost je velika, znaci ne postoji znacajna razlika izmedju dva modela
85 summary(glm(REMISS ~ LI, leukemia, family=binomial))
86 # AIC za ovaj model je manji nego za prvobitni, dakle bolji je novi model.

```