

Lineární modely včetně zobecněných, logistická regrese, analýza přežívání v R

—
17VSADR – Skriptování a analýza dat v jazyce R

Lubomír Štěpánek^{1, 2}



¹Oddělení biomedicínské statistiky
Ústav biofyziky a informatiky
1. lékařská fakulta
Univerzita Karlova v Praze



²Katedra biomedicínské informatiky
Fakulta biomedicínského inženýrství
České vysoké učení technické v Praze

(2019) Lubomír Štěpánek, CC BY-NC-ND 3.0 (CZ)



Dílo lze dále svobodně šířit, ovšem s uvedením původního autora a s uvedením původní licence. Dílo není možné šířit komerčně ani s ním jakkoliv jinak nakládat pro účely komerčního zisku. Dílo nesmí být jakkoliv upravováno. Autor neručí za správnost informací uvedených kdekoli v předložené práci, přesto vynaložil nezanedbatelné úsilí, aby byla uvedená fakta správná a aktuální, a práci sepsal podle svého nejlepšího vědomí a svých „nejlepších“ znalostí problematiky.

Obsah

Zavedení lineární regrese

- buď $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$ vektor spojité závisle proměnné, $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_k)^T$ je vektor lineárních koeficientů, nakonec \mathbf{X} je designová matice (též datová matice či matice modelu) právě k nezávisle proměnných (spojitých či kategorických) tak, že

$$\mathbf{X} = \begin{pmatrix} 1 & x_{1,1} & x_{1,2} & \cdots & x_{1,k} \\ 1 & x_{2,1} & x_{2,2} & \cdots & x_{2,k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & x_{n,2} & \cdots & x_{n,k} \end{pmatrix}$$

- pak lineární regresí (vícerozměrnou pro $k > 1$) nazveme model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

kde $\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)^T$ je vektor chybové složky (též méně přesně vektor reziduí)

Hledání řešení modelu lineární regrese

- iniciálně známe vektor \mathbf{y} a matici \mathbf{X}
- hledáme takový vektor odhadů lineárních koeficientů $\hat{\boldsymbol{\beta}} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k)^T$, aby součet druhých mocnin členů vektoru chybové složky $\boldsymbol{\varepsilon}$ byl co nejmenší
- tedy nalézt model lineární regrese mezi danými k nezávisle proměnnými a jednou závisle proměnnou znamená nalézt odhady lineárních koeficientů $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ tak, aby

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^n \varepsilon_i^2 \right\} = \arg \min_{\boldsymbol{\beta}} \{ \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} \} = \\ &= \arg \min_{\boldsymbol{\beta}} \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \} = \\ &= \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \dots - \beta_k x_{ik})^2 \right\}\end{aligned}$$

Tvar řešení lineární regrese

- řešení vede na soustavu normalizovaných rovnic, jejichž analytickým (!) řešením je vektor odhadů lineárních koeficientů

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- za předpokladu, že vektor chybové složky $\varepsilon = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)^T$ splňuje tzv. slabou sadu předpokladů, tj.
 - (i) nulová střední hodnota každého rezidua, tj. $\mathbf{E}(\varepsilon_i) = 0$ pro $\forall i \in \{1, 2, \dots, n\}$
 - (ii) (homoskedasticita) konečný konstantní rozptyl každého rezidua, tj. $\text{var}(\varepsilon_i) = \sigma^2 < \infty$ pro $\forall i \in \{1, 2, \dots, n\}$
 - (iii) (nekorelovanost) lineární nezávislost dvou různých reziduí, tj. $\text{cov}(\varepsilon_i, \varepsilon_j) = 0$ pro $\forall i, j \in \{1, 2, \dots, n\}$ a $i \neq j$
- lze ukázat, že řešení $\hat{\beta}$ má kýžené vlastnosti, např.

$$\mathbf{E}(\hat{\beta}_j - \beta_j) = 0$$

pro $\forall j \in \{0, 1, 2, \dots, k\}$ apod.

Tvar řešení lineární regrese

- současně lze nahlédnout, že odhady tzv. vyrovnaných hodnot \mathbf{y} vyjádříme za slabé sady předpokladů jako

$$\hat{\mathbf{y}} = \mathbf{E}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{E}(\mathbf{X}\boldsymbol{\beta}) + \mathbf{E}(\boldsymbol{\varepsilon}) = \mathbf{E}(\mathbf{X}\boldsymbol{\beta}) + \mathbf{0} = \mathbf{X}\mathbf{E}(\boldsymbol{\beta}) = \mathbf{X}\hat{\boldsymbol{\beta}}$$

a po rozepsání

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y},$$

kde tvar $\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ nazýváme též projekční či hat maticí a značíme $\mathbf{H} \equiv \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$

- znalost projekční matice \mathbf{H} je výhodná pro zkoumání vlivu závisle proměnné mezi různými pozorováními
- vidíme tedy, že jak lineární koeficienty $\hat{\boldsymbol{\beta}}$, tak vyrovnané hodnoty $\hat{\mathbf{y}}$ odhadneme pouze s pomocí apriorně známých hodnot, tj. datové matice \mathbf{X} a vektoru závisle proměnné \mathbf{y}

Některé důsledky pro výpočetní statistiku a analýzu dat

- protože $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ a $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$, všechny výpočty jsou maticové, tj. nevyžadují striktně numerické přístupy a obvykle existuje jednoznačné číselné řešení
- maticové výpočty mohou být náročné na paměť, méně na výpočetní čas – to je rozdíl oproti zobecněným lineárním či nelineárním modelům (time-memory trade-off)
- caveat může nastat ve fázi výpočtu $(\mathbf{X}^T \mathbf{X})^{-1}$, tedy inverze k $(\mathbf{X}^T \mathbf{X})$; má-li $(\mathbf{X}^T \mathbf{X})$ nízkou hodnotu v důsledku např. multikolinearity či *numerical fuzz*, je řešení číselně „nestabilní“ nebo vůbec neexistuje

Děkuji za pozornost!

lubomir.stepanek@lf1.cuni.cz

lubomir.stepanek@fbmi.cvut.cz

► GitHub

github.com/LStepanek/17VSADR_Skriptovani_a_analyza_dat_v_jazyce_R