

0 Lineární regrese

V modelu lineární regrese předpokládáme lineární závislost vysvětlované proměnné na hodnotách příznaků.

0.1 Model

Pro data s hodnotami příznaků X_1, \dots, X_p a hodnotou vysvětlované proměnné Y předpokládáme lineární model

$$Y = w_0 + w_1X_1 + \dots + w_pX_p + \varepsilon$$

kde w_i jsou neznámé koeficienty a ε představuje chybu nebo nekonzistenci výsledné hodnoty Y , kterou model nezachycuje a pro kterou platí $E\varepsilon = 0$.

V bodě $(x_1, \dots, x_p)^T$ tohoto modelu platí vztah

$$Y = w_0 + w_1x_1 + \dots + w_px_p + \varepsilon = \mathbf{w}^T \mathbf{x} + \varepsilon$$

kde zavádíme následující vektorovou notaci pro vstup \mathbf{x} a vektor vah \mathbf{w}

$$\begin{aligned}\mathbf{x} &= (1, x_1, \dots, x_p)^T \\ \mathbf{w} &= (w_0, w_1, \dots, w_p)^T\end{aligned}$$

0.2 Predikce

S odhadnutými váhami $\hat{\mathbf{w}}$ predikujeme vztahem

$$\hat{Y} = \hat{\mathbf{w}}^T \mathbf{x} = \hat{w}_0 + \hat{w}_1x_1 + \dots + \hat{w}_px_p$$

Pro skutečnou hodnotu $Y = \mathbf{w}^T \mathbf{x} + \varepsilon$ platí z předpokladu $E\varepsilon = 0$

$$EY = E\mathbf{w}^T \mathbf{x} + E\varepsilon = \mathbf{w}^T \mathbf{x}$$

\hat{Y} je tedy bodovým odhadem EY v bodě \mathbf{x} .

0.3 Metoda nejmenších čtverců

Metoda nejmenších čtverců nalézá hodnotu $\hat{\mathbf{w}}$ odhadu \mathbf{w} minimalizací následující kvadratické ztrátové funkce:

$$L(Y, \hat{Y}) = (Y - \hat{Y})^2$$

Pro trénovací množinou $(\mathbf{x}_i, Y_i), i = 1, \dots, N$ minimalizujeme reziduální součet čtverců

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N L(Y_i, \hat{Y}_i) = \sum_{i=1}^N (Y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

0.3.1 Maticový zápis trénovací množiny

Zavádíme náhodné vektory $\mathbf{Y} = (Y_1, \dots, Y_N)^T$, $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)$ a body $\mathbf{x}_1, \dots, \mathbf{x}_N$ zapisujeme do řádků matice

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} = \begin{pmatrix} 1 & x_{1,1} & \cdots & x_{1,p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & \cdots & x_{N,p} \end{pmatrix} \in \mathbb{R}^{N,p+1}$$

Při tomto značení platí rovnice $\mathbf{Y} = \mathbf{X}\mathbf{w} + \boldsymbol{\varepsilon}$ kde $E\boldsymbol{\varepsilon} = \mathbf{0}$.

0.3.2 Minimalizace RSS

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N (Y_i - \mathbf{w}^T \mathbf{x}_i)^2 = \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|^2$$

Začneme nalezením gradientu této funkce

$$\begin{aligned} \frac{\partial \text{RSS}(\mathbf{w})}{\partial w_j} &= \sum_{i=1}^N 2(Y_i - \mathbf{w}^T \mathbf{x}_i)(-x_{i,j}) \\ \nabla \text{RSS}(\mathbf{w}) &= \sum_{i=1}^N 2(Y_i - \mathbf{w}^T \mathbf{x}_i)(-\mathbf{x}_i) \\ &= -2\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\mathbf{w}) \end{aligned}$$

Položením $\nabla \text{RSS}(\mathbf{w}) = \mathbf{0}$ obdržíme normální rovnici

$$\mathbf{X}^T \mathbf{Y} - \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{0}.$$

Předpokládáme-li, že $\mathbf{X}^T \mathbf{X}$ je regulární, pak lze \mathbf{w} odhadnout následovně

$$\hat{\mathbf{w}}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

V případě regulární matice $\mathbf{X}^T \mathbf{X}$ je $\hat{\mathbf{w}}_{\text{OLS}}$ jediným kritickým bodem. Z geometrické interpretace, kterou dosáhneme stejného výsledku, bude plynout, že se jedná o globální minimum $\text{RSS}(\mathbf{w})$.

Predikce \hat{Y} v bodě \mathbf{x} je $\hat{Y} = \hat{\mathbf{w}}_{\text{OLS}}^T \mathbf{x}$.

0.3.3 Geometrická interpretace

Minimalizace $\text{RSS}(\mathbf{w}) = \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|^2$ je problém ekvivalentní minimalizaci $\|\mathbf{Y} - \mathbf{X}\mathbf{w}\|$. Hledáme tedy $\mathbf{w} \in \mathbb{R}^{p+1}$ takové, že vektory \mathbf{Y} a $\mathbf{X}\mathbf{w}$ jsou si v prostoru \mathbb{R}^N co nejblíže.

$\mathbf{X}\mathbf{w}$ je lineární kombinací sloupců \mathbf{X} :

$$\mathbf{X}\mathbf{w} = \sum_{i=0}^p \mathbf{w}_i \mathbf{X}_{:,i} \in \langle \mathbf{X}_{:,0}, \mathbf{X}_{:,1}, \dots, \mathbf{X}_{:,p} \rangle = \mathbf{P}$$

Hledaný bod $\mathbf{X}\mathbf{w}$ je proto nejbližší k bodu \mathbf{Y} , právě pokud je vektor $\mathbf{Y} - \mathbf{X}\mathbf{w}$ ortogonální na podprostor \mathbf{P} :

$$\mathbf{X}_{:,i}^T (\mathbf{Y} - \mathbf{X}\mathbf{w}) = 0 \text{ pro všechna } i = 0, 1, \dots, p$$

Což po přepsání do maticového tvaru dá opět normální rovnici

$$\mathbf{X}^T (\mathbf{Y} - \mathbf{X}\mathbf{w}) = \mathbf{X}^T \mathbf{Y} - \mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{0}$$

Z úvah také navíc plyne, že každé $\hat{\mathbf{w}}$, které splňuje normální rovnici, $\text{RSS}(\mathbf{w})$ minimalizuje (jedná se o globální minimum).

0.4 Regularita versus lineární nezávislost sloupců matice \mathbf{X}

Pro libovolné $\mathbf{s} \in \mathbb{R}^{p+1}$ platí následující posloupnost implikací:

$$\mathbf{X}^T \mathbf{X}\mathbf{s} = \mathbf{0} \Rightarrow \mathbf{s}^T \mathbf{X}^T \mathbf{X}\mathbf{s} = \|\mathbf{X}\mathbf{s}\|^2 = 0 \Rightarrow \mathbf{X}\mathbf{s} = \mathbf{0} \Rightarrow \mathbf{X}^T \mathbf{X}\mathbf{s} = \mathbf{0}$$

Z toho vyplývá, že $\mathbf{X}^T \mathbf{X}$ je regulární právě tehdy, když $\mathbf{X}\mathbf{s} = \mathbf{0}$ pouze pro $\mathbf{s} = \mathbf{0}$, což platí právě tehdy, když jsou sloupce matice \mathbf{X} lineárně nezávislé. To zřejmě nemusí vždy platit: např. když $N < p + 1$ nebo pokud je nějaký příznak lineární kombinací ostatních.

Normální rovnice má v případě regulární $\mathbf{X}^T \mathbf{X}$ právě jedno řešení.

$$\hat{\mathbf{w}}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

V opačném případě má jich má nekonečně mnoho.

0.5 Problém kolinearity

Pokud je nějaký příznak “skoro” lineární kombinací ostatních příznaků, pak můžeme narazit na problém kolinearity. V případě silně korelovaného příznaku je příslušná proměnná do jisté míry zbytečná, protože ve svém rozměru (dimenzi) nepřidává žádnou informaci navíc a zároveň může mít ve výsledném vektoru $\hat{\mathbf{w}}_{\text{OLS}}$ vysoký koeficient. Jako příklad lze uvést 3-dimenzionální prostor (X_1, X_2, Y) ve kterém řešení degeneruje do skoro-přímky, kterou lze protnou rovinou více způsoby.

0.5.1 Důsledky kolinearity

Ve výsledku kolinearita způsobuje vysoký rozptyl $\hat{\mathbf{w}}_{\text{OLS}}$, který je tak velmi citlivý na data. Pro různé realizace stejného náhodného výběru se řešení normální rovnice může značně lišit. Vysoký rozptyl se následovně přenáší na predikce, které jsou pak méně spolehlivé.

0.5.2 Řešení problému kolinearity

Řešením by bylo odstranit příznaky, které kolinearity způsobují. To však není vždy jednoduché. S vyšším počtem příznaků je velmi obtížné takové sloupce identifikovat. Také se může stát, že kolinearita je mezi velkým počtem příznaků, které jsou dohromady klíčové pro správnou predikci. V takovém případě není ideální takové příznaky zahodit. Existují však metody, které dokážou počet příznaků snížit (odstraněním, případně nahrazením menším počtem) tak, aby byly sloupce LN.

Je také možné změnit funkci, kterou minimalizujeme, abychom měli stabilnější a jednoznačnější řešení. Typicky se přidává regulační člen, který problémy kolinearity může zmírnit (hřebenová regrese v sekci ?? nebo lasso v sekci ??).

0.6 Shrnutí

Lineární regrese je rezistentní vůči problémům spojené s vysokou dimenzí dat. Problém nastává, když je dat méně než příznaků nebo když jsou příznaky silně korelované.