



درس نامه چهارم

۱ رگرسیون چند بعدی

در مساله رگرسیون چندبعدی، ورودی یک بردار p -بعدی و خروجی یک عدد می باشد.



به بعدهای ورودی:

- Input
- Predictor
- Feature
- Variable Independent

و به بعد خروجی اصطلاحات:

- Output
- Response
- Variable Dependent

اطلاق میگردد.

معمولا داده را به صورت یک ماتریس نشان میدهیم:

x_1	x_2	...	x_p	y
		...		

مدل به صورت زیر در نظر گرفته میشود:

$$\hat{y} = \beta_0 + \sum_{j=1}^p \beta_j x_j \quad (1)$$

تابع ریسک (یا خسارت) را نیز به صورت:

$$l(y, \hat{y}) = (y - \hat{y})^2 \quad (2)$$

تعریف میکنیم. بنابراین برای کمینه کردن ریسک (یا خسارت) از تعریف زیر استفاده میکنیم:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \underset{\beta}{\operatorname{argmin}} \mathbf{RSS}(\beta) \quad (3)$$

که در آن داریم:

$$\beta = (\beta_0, \beta_1, \dots, \beta_p)^T \quad (4)$$

و

$$\mathbf{RSS}(\beta) = \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2. \quad (5)$$

رابطه فوق را میتوان به صورت ماتریسی نیز نشان داد. اگر ماتریس X را به صورت زیر تعریف کنیم:

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{np} \end{bmatrix} \quad (6)$$

و بردار y را به صورت زیر تعریف کنیم:

$$y = (y_1, \dots, y_n)^T \quad (7)$$

میتوانیم عبارت RSS را به صورت زیر بازنویسی کنیم:

$$\mathbf{RSS}(\beta) = (y - X\beta)^T (y - X\beta) = \|y - X\beta\|_2^2 \quad (8)$$

در نتیجه اگر نسبت به $\hat{\beta}$ مشتق بگیریم و در $\hat{\beta}$ برابر صفر قرار دهیم:

$$\frac{\partial RSS}{\partial \beta} \Big|_{\beta=\hat{\beta}} = -2X^T(y - X\hat{\beta}) = 0 \quad (9)$$

به دستگاه خطی زیر می‌رسیم:

$$X^T X \hat{\beta} = X^T y \quad (10)$$

اگر مشتق دوم RSS را محاسبه کنیم:

$$\frac{\partial^2 RSS}{\partial \beta \partial \beta^T} = 2X^T X \quad (11)$$

خواهیم دید مشتق دوم، به ازای تمام β -ها یک ماتریس مثبت نیمه معین است. (چرا؟) بنابراین تابع $RSS(\beta)$ یک تابع محدب است و جواب‌های ۱۰ همگی نقاط بهینه می‌باشند.

اگر X دارای رتبه ستونی کامل باشد، آنگاه $X^T X$ رتبه کامل (Full Rank) خواهد بود. پس $\hat{\beta}$ یک جواب یکتا دارد.

در این حالت:

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad (12)$$

و مقدار \hat{y} به صورت زیر قابل محاسبه است.

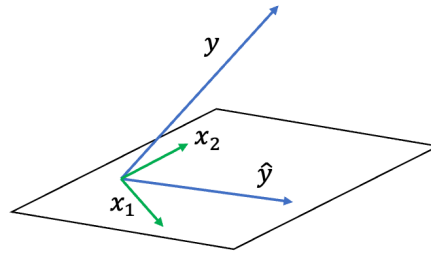
$$\hat{y} = X\hat{\beta} = \underbrace{X(X^T X)^{-1} X^T}_H y \quad (13)$$

۲ دیدگاه هندسی

در حقیقت مساله از دیدگاه هندسی بدین صورت است که می‌خواهیم یک بردار \hat{y} را که از ترکیب خطی ستونهای X حاصل میشود، طوری بیابیم که کمترین فاصله را با y داشته باشد. این بدین معنی است که باید y را روی فضایی که با ستونهای X ساخته میشود، تصویر نماییم:

$$\hat{y} = P_{col(1, x_1, \dots, x_p)}(y) \quad (14)$$

ماتریس H در حقیقت این عمل را انجام میدهد.



در شرایطی ممکن است ستونهای ماتریس X مستقل خطی نباشد. در این شرایط هرچند \hat{y} مقداری مشخص است اما برای $\hat{\beta}$ جوابهای متعددی وجود خواهد داشت.

یکی از راههای مقابله با این مساله این است که ستونهایی که به کمک ستونهای دیگر قابل بدست آمدن میباشد (ترکیب خطی ستونهای دیگر است) را دور بریزیم تا ماتریس X دارای رتبه ستونی کامل شود.

۳ ارزیابی

اگر m نمونه را کنار گذاشته باشیم، میتوانیم تعمیم را مقایسه ریسک (یا خسارت) درون داده و بیرون داده مورد آزمون قرار دهیم.

اگر مدل واقعی به صورت زیر باشد:

$$y = \underbrace{(1, x_1, \dots, x_p)}_{X^T} \beta^* + \epsilon \quad (15)$$

خواهیم داشت:

$$\mathbb{E}[l(y, \hat{y})] = \mathbb{E}[(X^T(\beta^* - \hat{\beta}) + \epsilon)^2] \quad (16)$$

$$= \mathbb{E}[\{(X^T(\beta^* - \mathbb{E}[\hat{\beta}]) + X^T(\mathbb{E}[\hat{\beta}] - \hat{\beta}) + \epsilon\}^2] \quad (17)$$

$$= \|X^T(\beta^* - \mathbb{E}[\hat{\beta}])\|_2^2 + \mathbb{E}[\|X^T(\mathbb{E}[\hat{\beta}] - \hat{\beta})\|_2^2] + \sigma_\epsilon^2 \quad (18)$$

$$= \text{Bias}^2 + \text{Variance} + \text{Noise} \quad (19)$$

به این صورت میتوان ریسک (خسارت) را به سه جز مذکور تجزیه کرد. (Bias-Variance Decomposition)

۴ استنتاج

اگر $\epsilon_i \sim \mathcal{N}(0, \sigma_\epsilon^2)$ را فرض کنیم، خواهیم داشت:

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad (20)$$

$$= (X^T X)^{-1} X^T (X\beta^* + \epsilon) \quad (21)$$

$$= \beta^* + (X^T X)^{-1} X^T \epsilon \quad (22)$$

بنابراین با فرض دانستن X داریم:

$$\hat{\beta} \sim N(\beta^*, (X^T X)^{-1} \sigma_\epsilon^2) \quad (23)$$

بنابراین، $\hat{\beta}$ دارای توزیع گوسی میباشد. تخمین σ_ϵ^2 را نیز میتوان به صورت زیر محاسبه کرد:

$$\hat{\sigma}_\epsilon^2 = \frac{\text{RSS}}{n - p - 1} = \frac{1}{n - p - 1} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (24)$$