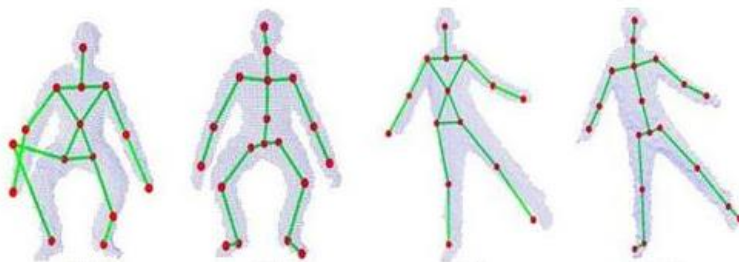


بسمه تعالی



دانشگاه صنعتی شریف
دانشکده مهندسی برق
گروه سیستم‌های دیجیتال



آزمایشگاه یادگیری و بینایی ماشین

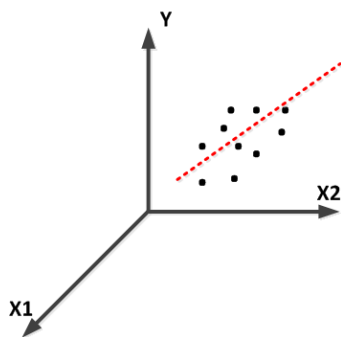
دستور کار آزمایش دوم: رگرسیون خطی

زمان لازم برای انجام آزمایش: حداکثر دو جلسه

Sajjad Hashembeiki
98107077

قسمت اول: رگرسیون خطی با تعداد ویژگی های محدود

منظور از رگرسیون خطی این است که با دانستن نقاطی نمونه از داده آماری، بتوان رفتار آن داده را در سایر نقاط و یا در آینده پیش بینی کرد. برای مثال، فرض کنید که کمیت (y) ما آلودگی هوا است. این آلودگی هوا به عواملی از جمله میزان تردد (x_1) و وسعت فضای سبز (x_2) وابسته است. در این صورت، اگر از تعدادی از نقاط شهر، نمونه ای از آلودگی هوا را بگیریم، می توانیم نمودار زیر را تشکیل دهیم.



همان طور که دیده می شود، تا وقتی که داده جمع آوری شده، کاملاً نویز نباشد، می توان یک الگوی مشخصی از آن بدست آورد و یک خط با شیب غیر صفر به آن برازش کرد.

این خط برازش شده را در شکل، به صورت یک خط چین نمایش داده ایم.

اگر متغیرهای مساله را به صورت یک بردار x نمایش دهیم، خواهیم داشت:

$$x = [x_1, x_2, \dots, x_p]^T$$

در این صورت، خروجی ما (که از آن نمونه گرفته ایم) به صورت زیر است:

$$y = f(x) + \epsilon$$

که در آن، ϵ نویز است. هدف ما یافتن $f(\cdot)$ است که همان رگرسیون خطی داده است:

$$f(x) = \beta_0 + \sum_{j=1}^p \beta_j x_j$$

فرض کنید که N نمونه گرفته ایم یا داریم. که متغیرهای نمونه i ام به صورت زیر است:

$$x^i = [x_1^i, x_2^i, \dots, x_p^i]^T$$

و خروجی آن (نمونه i ام) به صورت y^i است.

برای یافتن f ما باید ضرایب β را پیدا کنیم. باید β هایی را بیابیم که کمترین اختلاف بین y^i ها و $f(x^i)$ ها را ایجاد کند تا بهترین تخمین ممکن را داشته باشیم. ما برای نمایش تفاضل این دو مقدار، از تفاضل نرم ۲ استفاده می‌کنیم:

$$RSS = \sum_{i=1}^N (y^i - f(x^i))^2$$

ما باید عبارت بالا را حداقل کنیم. می‌توان نشان داد که عبارت RSS بالا به بیان ماتریسی، به صورت زیر می‌شود:

$$RSS = (y - X\beta)^T (y - X\beta)$$

که در آن، ماتریس X به صورت زیر است:

$$X = \begin{bmatrix} 1 & x^{1T} \\ \vdots & \vdots \\ 1 & x^{NT} \end{bmatrix}_{N \times (P+1)}$$

و y بردار ستونی است که خروجی نمونه‌ها را نشان می‌دهد.

(امتیازی): رابطه ماتریسی RSS را اثبات نمایید.



اگر از رابطه RSS بالا نسبت به β مشتق بگیریم و مساوی صفر قرار دهیم (مساله بهینه‌سازی)، به صورت زیر بدست می‌آید:

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

از رابطه بالا برای تخمین y (خروجی نمونه‌های آموزشی) به صورت زیر استفاده می‌شود:

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

Bonus

$$RSS = \sum_{i=1}^n (y^i - f(x^i))^2$$

$$f(x) = \beta_0 + \sum_{j=1}^p \beta_j x_j$$

$$\underline{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}_{(p+1) \times 1}, \quad \underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}_{p \times 1}, \quad \underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

$$\Rightarrow f(x) = \underline{x}^T \underline{\beta}$$

$$RSS = \sum_{i=1}^n (y^i - \underline{x}^{iT} \underline{\beta})^2$$

$$= \begin{bmatrix} y^1 - \underline{x}^{1T} \underline{\beta} & y^2 - \underline{x}^{2T} \underline{\beta} & \dots & y^n - \underline{x}^{nT} \underline{\beta} \end{bmatrix}_{1 \times n} \begin{bmatrix} y^1 - \underline{x}^{1T} \underline{\beta} \\ y^2 - \underline{x}^{2T} \underline{\beta} \\ \vdots \\ y^n - \underline{x}^{nT} \underline{\beta} \end{bmatrix}_{n \times 1} = Z^T Z$$

$$\underline{Z} = \begin{bmatrix} y^1 - \underline{x}_1^T \underline{\beta} \\ y^2 - \underline{x}_2^T \underline{\beta} \\ \vdots \\ y^n - \underline{x}_n^T \underline{\beta} \end{bmatrix}_{n \times 1} = \underline{y} - \begin{bmatrix} \underline{x}_1^T \underline{\beta} \\ \underline{x}_2^T \underline{\beta} \\ \vdots \\ \underline{x}_n^T \underline{\beta} \end{bmatrix}_{n \times 1}$$

$$= \underline{y} - \begin{bmatrix} 1 & x_1^1 & x_2^1 & \dots & x_p^1 \\ 1 & x_1^2 & x_2^2 & \dots & x_p^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^n & x_2^n & \dots & x_p^n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}_{(p+1) \times 1}$$

$n \times (p+1) \quad (p+1) \times 1$

$$= \underline{y} - \underline{X} \underline{\beta}$$

$$\Rightarrow RSS = \underline{Z}^T \underline{Z} \xrightarrow{\underline{Z} = \underline{y} - \underline{X} \underline{\beta}} \boxed{RSS = (\underline{y} - \underline{X} \underline{\beta})^T (\underline{y} - \underline{X} \underline{\beta})}$$

پیش گزارش

۱- در یک مسأله رگرسیون خطی، RSS را به صورت زیر تعریف می شود:

$$RSS = \sum_{i=1}^n (y^i - \beta_0 - \beta_1 x^i)^2$$

β_0^* و β_1^* جواب‌های این مسأله هستند. از بین معادله‌های زیر، آن‌هایی که صحیح هستند را مشخص کنید.
(راهنمایی: از معادله بالا نسبت به β_0 و β_1 مشتق بگیرید.)

$$\begin{array}{l} \text{False} \left\{ \begin{array}{l} \sum_{i=1}^n (y^i - \beta_0^* - \beta_1^* x^i) y^i = 0 \quad \times \\ \sum_{i=1}^n (y^i - \beta_0^* - \beta_1^* x^i) (y^i - \bar{y}) = 0 \quad \times \end{array} \right. \\ \text{True} \left\{ \begin{array}{l} \sum_{i=1}^n (y^i - \beta_0^* - \beta_1^* x^i) (x^i - \bar{x}) = 0 \quad \checkmark \\ \sum_{i=1}^n (y^i - \beta_0^* - \beta_1^* x^i) (\beta_0^* + \beta_1^* x^i) = 0 \quad \checkmark \end{array} \right. \end{array}$$

۲- در دستور کار با مفهوم **overfitting** آشنا خواهید شد. در مورد **underfitting** تحقیق نمایید و کمی درباره این مفهوم توضیح دهید.

$$1) \text{RSS} = \sum_{i=1}^n (y^i - \beta_0 - \beta_1 x^i)^2$$

$$\rightarrow \frac{\partial \text{RSS}}{\partial \beta_0} = \sum_{i=1}^n -1 \times 2 (y^i - \beta_0 - \beta_1 x^i)$$

$$\Rightarrow \underbrace{\sum_{i=1}^n (y^i - \beta_0^* - \beta_1^* x^i)}_{=A} = 0 \quad (*)$$

$$\rightarrow \frac{\partial \text{RSS}}{\partial \beta_1} = \sum_{i=1}^n -x^i \times 2 (y^i - \beta_0 - \beta_1 x^i)$$

$$\Rightarrow \underbrace{\sum_{i=1}^n x^i (y^i - \beta_0^* - \beta_1^* x^i)}_{=B} = 0 \quad (**)$$

The third eqn $\Rightarrow B - \bar{x} A \xrightarrow{B=0, A=0} 0 \checkmark$

The last eqn $\Rightarrow A \beta_0^* + \beta_1^* B \xrightarrow[A=0]{B=0} 0 \checkmark$

2: Underfitting: model can't capture (learn) the relationship between inputs & outputs properly and error rate on training set & test set is high. Basically the model isn't complex enough. In this situation the model has High Bias and low variance.

