



## یادگیری ماشین

نیم‌سال دوم ۱۴۰۱-۱۴۰۲

مدرس: دکتر سید ابوالفضل مطهری

## دورسی فاصله پنجم

بنابراین  $\hat{\beta}$  دارای یک توزیع گوسی است. تخمین  $\sigma_\epsilon^2$  را می‌توانیم از روی رابطه زیر داشته باشیم:

$$\hat{\sigma}_\epsilon^2 = \frac{1}{N - P - 1} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

که می‌توان نشان داد که توزیع  $\hat{\sigma}_\epsilon^2$  یک توزیع chi-squared می‌باشد یعنی:

$$(N - P - 1)\hat{\sigma}_\epsilon^2 \sim \sigma_\epsilon^2 \chi_{N-P-1}^2$$

همچنین می‌توان نشان داد که  $\hat{\sigma}_\epsilon^2$  و  $\hat{\beta}$  از دیدگاه آماری مستقل از یکدیگر می‌باشند.

## آزمون فرض

اگر بخواهیم  $\beta_j = 0$  را تست نماییم مقدار زیر را محاسبه می‌کنیم:

$$t_{n-p-1} = \frac{\hat{\beta}_j}{\hat{\sigma}_\epsilon \sqrt{v_j}}$$

مقدار فوق دارای توزیع  $t$  با درجه آزادی  $n - p - 1$  می‌باشد که در آن  $v_j$ ،  $j$ -امین مقدار قطری  $(X^T X)^{-1}$  می‌باشد. اگر بخواهیم به طور گروهی آزمون

$$\beta_j = \beta_k = \dots = 0$$

را مورد ارزیابی قرار دهیم آنگاه یکبار مدل را با مقادیر صفر فوق یاد می‌گیریم و سپس با تمامی مقدارها. در حالت اول اگر میزان مجموع مربعات باقیمانده را  $RSS_0$  و در حالت دوم  $RSS_1$  نمایش دهیم آنگاه:

$$F = \frac{(RSS_0 - RSS_1)/p_1 - p_0}{RSS_1/n - p_1 - 1}$$

دارای توزیع  $F$  با پارامترهای  $(p_1 - p_0, n - p_1 - 1)$  می‌باشد.

## متعامد سازی

با توجه به آنکه تجزیه  $X = QR$  از طریق گرم-اشمیت و دیگر روش ها قابل محاسبه است که در آن  $Q^T Q = I$  و  $R$  یک ماتریس بالا مثلثی است خواهیم داشت :

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

اگر ماتریس  $Q$  را به صورت

$$Q = [q_1, q_2, \dots, q_p]$$

و ماتریس  $R^{-1}$  را به صورت

$$R = \begin{bmatrix} \frac{r_{11}}{r_{11}} & & & \\ & r_{11} & & \\ & & \ddots & \\ & & & \end{bmatrix}$$

نمایش دهیم ، رابطه فوق نشان می دهد:

$$\begin{aligned} r_{11} \hat{\beta}_1 &= q_1^T y \\ r_{11} \hat{\beta}_1 + r_{12} \hat{\beta}_2 &= q_2^T y \\ r_{12} \hat{\beta}_1 + r_{22} \hat{\beta}_2 &= q_3^T y \\ &\vdots \end{aligned}$$

رابطه فوق را می توان به این صورت دید که رابطه اول در حقیقت یک رگرسیون ساده و یک بعدی برای داده های  $x_1 = (1, \dots, 1)$  و  $y = (y_1, \dots, y_n)$  می باشد .  $\hat{\beta}_1$  از روی این مسئله به دست می آید. سپس برای  $\hat{\beta}_2$  بایستی مسئله رگرسیون ساده دوم را حل نماییم و الی آخر.

**نکات مهم :**

۱. اگر یکی از بعد ها داده های دسته بند داشت می بایست آنرا کسر نماییم . به این منظور داده های دو تایی مانند داشتن حساب و نداشتن حساب را یا با  $(0, 1)$  و یا  $(-1, 1)$  کد می نماییم . اگر تعداد دسته ها بیشتر بود می توانیم چند ستون اضافه کنیم و به ازای وجود هر کدام از آن دسته ها عدد ۱ و در غیر این صورت صفر قرار می دهیم .

۲. ورودی های جدید را می توانیم بر پایه ورودی های اولیه تولید کنیم و بدین ترتیب پیچیدگی را بالا ببریم، در این شرایط هنوز به خاطر وجود ضرایب خطی با یک رگرسیون خطی مواجه هستیم. مثلاً:

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4 x_1^2 + \beta_5 x_1^3$$

## مشکلات احتمالی

۱. رابطه واقعی ممکن است غیر خطی باشد.

اگر رابطه واقعی غیر خطی باشد آنگاه

$$e = y - \hat{y} = f(x) - \langle x, \hat{\beta} \rangle + \varepsilon$$

بنابراین  $e$  و  $x$  مستقل نمی شوند و در این صورت اگر رابطه  $e$  و  $x$  را بکشیم یک رابطه را مشاهده می کنیم. اگر  $x$  چند بعدی بود آنگاه می توانیم از  $\hat{y} = \langle x, \hat{\beta} \rangle$  استفاده نماییم و رابطه آن را با  $e$  بررسی کنیم. به هر حال در حالت ایده آل  $e$  فقط با  $\varepsilon$  رابطه دارد.

۲. واریانس خطا ممکن است ثابت نشود.

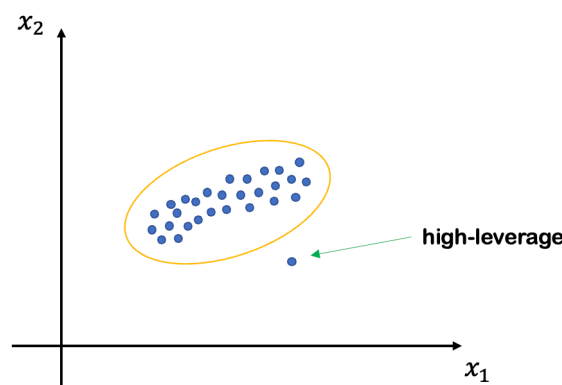
ممکن است فرض آنکه  $\varepsilon_i$  ها مستقل و یا ناهمبسته می باشند درست نباشد. در چنین شرایطی استنتاجها می تواند اشتباه باشند. مثلا در داده های سری زمانی ممکن است نقاط نزدیک دارای خطاهای اندازه گیری همبسته باشند. برای بررسی این قضیه می توان پس از رگرسیون خطا را برحسب زمان ترسیم نمود و اگر رابطه ای در خطا دیده شد آن را یافت. البته در مسائل غیر سری زمانی یافتن آنها مشکل است.

۳. وجود نمونه های ناهنجار

نمونه های ناهنجار، نمونه هایی هستند که از مدل نیامده اند و احتمالا به خطا وارد داده شده اند. برای از بین بردن آن می توانیم مقدار  $\frac{e_i}{\hat{SE}(e)}$  را برای تمامی  $i$  ها بدست آوریم. مقایسه مقدار فوق با یک آستانه می تواند نشان دهد که مقدار ناهنجار است یا خیر.

۴. نقاط High-leverage

به نقاطی اطلاق می گردد که ورودی  $x$  نامعقولی نسبت به بقیه نقاط داشته باشند.



این نقاط تاثیر زیادی بر یادگیری می گذارند و از بین بردن آنها می تواند تاثیر زیادی بر کیفیت داشته باشد. یکی از راه های یافتن آنها محاسبه آماره زیر است:

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i'=1}^n (x_{i'} - \bar{x})^2}$$

می‌دانیم  $1 \geq h_i \geq \frac{1}{n}$  بوده و  $\frac{\sum h_i}{n} = \frac{p+1}{n}$ . بنابراین اگر  $h_i$  با مقدار خیلی بزرگتر از  $\frac{p+1}{n}$  یافت شد، می‌توان آن را حذف نمود.

## ۵. Collinearity

ممکن است تعدادی از ستون‌های داده با هم رابطه داشته باشند. یعنی همبستگی خطی بین آن‌ها دیده شود. در این شرایط نمی‌توان تاثیر این دو ستون را از یکدیگر تفکیک نمود. اگر به ماتریس همبستگی نگاه کنیم، مقداری که مقدار بزرگ دارند می‌توانند بیانگر چنین ستون‌هایی باشند. اما اگر بخواهیم تاثیر چندگانه را بیابیم از VIF (variance inflation factor) استفاده می‌نماییم.

$$\text{VIF}(\hat{\beta}_j) = \frac{1}{1 - R_{x_j|x_j}^2}$$

یعنی از بقیه پارامترها  $x_j$  را تخمین بزنیم.