

x	0.31634	0.22809	1.387	0.1712
z1	-24.24981	31.51425	-0.769	0.4450
z2	-34.10414	34.74505	-0.982	0.3307
x:z1	0.08521	0.31483	0.271	0.7877
x:z2	0.18995	0.35103	0.541	0.5907

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.741 on 54 degrees of freedom

Multiple R-squared: 0.4591, Adjusted R-squared: 0.409

F-statistic: 9.166 on 5 and 54 DF, p-value: 2.329e-06

همان‌گونه که در خروجی بالا می‌توان دید، فرض صفر بودن تمامی ضریب‌های مدل رگرسیون برازش داده شده غیر از عرض از مبدأ در سطح خطای ۵٪ پذیرفته شده است. علت این مسأله وجود همخطی چندگانه‌ی شدید بین ضریب‌های مدل رگرسیون با اثرهای متقابل است.

برای دیدن مقدارهای VIF در این مدل، تابع `vif()` در بسته‌ی نرم‌افزاری `car` را به صورت زیر به کار می‌بریم.

```
> library(car)
```

```
> vif(fit)
```

x	z1	z2	x:z1	x:z2
2.869837	139.542183	169.620140	139.148762	166.479597

در خروجی بالا می‌بینید، به‌طور تقریبی همه‌ی مقدارهای VIF از ۱۰۰ بزرگ‌تر هستند، که این نشان‌دهنده‌ی همخطی چندگانه‌ی شدید در مدل برازش داده شده است.

## ۱۱.۴ تمرین‌ها

۱. بر اساس داده‌های  $(x_1, y_1), \dots, (x_n, y_n)$ ، اطلاعات زیر به دست آمده است

$$\sum_{i=1}^n y_i = 1100, \quad \sum_{i=1}^n x_i y_i = 61800, \quad (X'X)^{-1} = \frac{1}{34000} \begin{bmatrix} 28400 & -500 \\ -500 & 10 \end{bmatrix};$$

که در آن  $X' = \begin{bmatrix} 1 & \cdots & 1 \\ x_1 & \cdots & x_n \end{bmatrix}$  است. مقدار  $\hat{\beta}$  را به دست آورید.

۲. مقدارهای  $\beta_0$  و  $\beta_1$  را به دست آورید.

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 2 & 0 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} = \begin{pmatrix} 10 \\ 0 \\ 5 \end{pmatrix}$$

۳. در مدل رگرسیون خطی چندگانه به صورت زیر

$$y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_{p-1} x_{i(p-1)} + \epsilon_i, \quad i = 1, \dots, n$$

با فرض  $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$ ، نشان دهید به ازای هر  $j \neq j'$ ،  $\hat{\beta}_j$  و  $\hat{\beta}_{j'}$  از هم مستقل هستند، اگر و تنها اگر

$$\sum_{i=1}^n x_{ij} x_{ij'} = 0, \quad j \neq j' (= 1, \dots, p-1)$$

۴. در مدل رگرسیونی  $y = X_n \beta + \epsilon$  با  $\hat{y} = X \hat{\beta}$ ، نشان دهید  $\text{Var}(\hat{y}) = \sigma^2 H$  و در نتیجه

$$\sum_{i=1}^n \text{Var}(\hat{y}_i) = p\sigma^2$$

۵. برای مدل رگرسیون خطی چندگانه، با فرض همگنی واریانس خطاها، نشان دهید

$$\text{Var}(\bar{y}) \leq \text{Var}(\hat{y}_i) \leq \text{Var}(y_i), \quad i = 1, \dots, n$$

۶. میانگین  $m \geq 1$  مقدار مشاهده شده‌ی مستقل از متغیر پاسخ برای یک سطح از متغیر پیشگو

$$(x^*) \text{ یعنی } \bar{y}_m^* = \frac{1}{m} \sum_{j=1}^m y_{jx^*} \text{ را در نظر بگیرید.}$$

الف. متغیر تصادفی  $Z^* = \bar{y}_m^* - \hat{y}_x^*$  را در نظر بگیرید. نشان دهید

$$\text{Var}(Z^*) = \sigma^2 \left[ \frac{1}{m} + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_x^2} \right]$$

ب. یک بازه‌ی پیش‌بینی  $\alpha(1 - \alpha) 100\%$  با دم‌های برابر برای  $\bar{y}_m^*$  به دست آورید.

۷. نشان دهید در مدل رگرسیون خطی ساده، درایه‌های ماتریس  $H = X(X'X)^{-1}X'$  به صورت زیر هستند.

$$h_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{S_{xx}}, \quad i, j = 1, \dots, n$$

۸. فرض کنید می‌خواهیم برآورد ضریب‌های  $\beta$  در مدل رگرسیونی  $y = X\beta + \epsilon$  را با در نظر گرفتن  $R\beta = r$  به دست آوریم. با تشکیل تابع لاگرانژ برای حداقل کردن مجموع توان دوم خطا با در نظر گرفتن قید داده شده به صورت زیر

$$\psi(\beta, \lambda) = (y - X\beta)'(y - X\beta) - 2\lambda'(R\beta - r)$$

نشان دهید که برآورد کم‌ترین توان دوم مقید  $\beta$  برابر است با

$$\hat{\beta}_R = \hat{\beta} + (X'X)^{-1}R'[R(X'X)^{-1}R']^{-1}(r - R\hat{\beta})$$

۹. برای مدل زیر هر یک از فرض‌های خطی زیر را به فرم کلی  $R\beta = 0$  با مشخص کردن ماتریس  $R$  بنویسید

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5 + \epsilon$$

الف.  $H_0: \beta_1 = \beta_2 = \beta_4$

ب.  $H_0: \beta_1 = \beta_2, \beta_3 = \beta_4, \beta_5 = 0$

ج.  $H_0: \beta_1 - 2\beta_5 = 3\beta_4, \beta_1 + 2\beta_3 = 0$

۱۰. داده‌های جدول ۹.۴ حجم بازدمی ( $y$ ) ۱۹ بیمار مبتلا به آسم را به همراه سن ( $x_1$ )، جنسیت ( $x_2$ )، قد ( $x_3$ ) و وزن ( $x_4$ ) آن‌ها ارائه می‌دهد (بتل و همکاران، ۱۹۸۵).



جدول ۹.۴: داده‌های بیماران مبتلا به آسم ( $0 = \text{زن}$ ,  $1 = \text{مرد}$ )

ردیف	حجم بازدمی ( $y$ )	سن ( $x_1$ )	جنسیت ( $x_2$ )	قد ( $x_3$ )	وزن ( $x_4$ )
۱	۴/۷	۲۴	۱	۱۷۵	۷۸/۰
۲	۴/۳	۳۶	۱	۱۷۳	۶۷/۶
۳	۳/۵	۲۸	۰	۱۷۱	۹۸/۰
۴	۴/۰	۲۵	۱	۱۶۶	۶۵/۵
۵	۳/۲	۲۶	۰	۱۶۶	۶۵/۰
۶	۴/۷	۲۲	۱	۱۷۶	۶۵/۵
۷	۴/۳	۲۷	۱	۱۸۵	۸۵/۵
۸	۴/۷	۲۷	۱	۱۷۱	۷۶/۳
۹	۵/۲	۳۶	۱	۱۸۵	۷۹/۰
۱۰	۴/۲	۲۴	۱	۱۸۲	۸۸/۲
۱۱	۳/۵	۲۶	۱	۱۸۰	۷۰/۵
۱۲	۳/۲	۲۹	۱	۱۶۳	۷۵/۰
۱۳	۲/۶	۳۳	۰	۱۸۰	۶۸/۰
۱۴	۲/۰	۳۱	۱	۱۸۰	۶۵/۰
۱۵	۴/۰	۳۰	۱	۱۸۰	۷۰/۴
۱۶	۳/۹	۲۲	۱	۱۶۸	۶۳/۰
۱۷	۳/۰	۲۷	۱	۱۶۸	۹۱/۲
۱۸	۴/۵	۴۶	۱	۱۷۸	۶۷/۵
۱۹	۲/۴	۳۶	۱	۱۷۳	۶۲/۰

الف. یک مدل رگرسیون خطی چندگانه به این داده‌ها برازش دهید.

ب. آزمون فرض معنی‌داری ضریب‌های تک‌تک متغیرهای پیشگو را انجام دهید.

ج. مقدار  $R^2$  و  $R^2_{adj}$  را محاسبه کنید و درباره‌ی نیکویی برازش مدل به داده‌ها اظهارنظر کنید.

۱۱. برای داده‌های تمرین ۱۰:

الف. عامل‌های تورم واریانس (VIF) ضریب‌های رگرسیونی را حساب کنید و در مورد همبستگی

خطی متغیرهای پیشگوی  $x_1, \dots, x_4$  اظهارنظر کنید.

ب. نمودارهای متغیرهای اضافه شده را رسم کنید و درباره‌ی خطی بودن مدل رگرسیون خطی

چندگانه اظهارنظر کنید.

ج. نمودارهای مانده‌ها را رسم کنید و درباره‌ی درستی فرض‌های اساسی اظهارنظر کنید.

۱۲. فرض کنید  $h_{ii}$  برابر  $i$ امین عضو قطر اصلی ماتریس  $H$  باشد. ثابت کنید در یک رگرسیون خطی چندگانه،  $\frac{1}{n} \leq h_{ii} \leq 1$  (از خودتوان بودن ماتریس  $H$  استفاده کنید).

۱۳. می‌دانیم که برای یک بردار تصادفی  $y$  (با  $n$  مؤلفه) با میانگین  $\mu$  و ماتریس واریانس-کواریانس  $\sigma^2 A$ ، که در آن  $A$  یک ماتریس معین مثبت  $n \times n$  است،

$$E(y' Ay) = \sigma^2 \text{tr}(A) + \mu' A \mu,$$

الف. نشان دهید

$$E(S_e^2) = \sigma^2$$

ب. در مدل رگرسیونی  $y = x\beta + \epsilon$  با  $\text{Cov}(\epsilon) = \sigma^2 I_n$  نشان دهید

$$E(y' y) = n\sigma^2 + \beta' X' X \beta.$$

۱۴. برای پیش‌بینی برخی مشخصه‌های گونه‌های خاص از درختان یک جنگل مانند تعداد و یا متوسط قطر آن درختان می‌توان برخی مشخصه‌های دیگر آن‌ها را به‌کارگرفت. در جدول ۱۰.۴ برخی مشخصه‌های درختان کاج در یک ناحیه جنگلی مانند سن درختان (Age)، متوسط ارتفاع درختان شاخص (HD)، تعداد درختان در یک جریب (N) و متوسط قطر درختان (MDBH) داده شده است (مایرز، ۱۹۹۰). این داده‌ها برای پیش‌بینی MDBH به عنوان متغیر پاسخ به کار گرفته می‌شوند و متغیرهای مستقل در آن عبارتند از  $x_1 = \text{HD}$ ،  $x_2 = \text{Age} \times N$  و  $x_3 = \frac{\text{HD}}{N}$ .

الف. مدل رگرسیون خطی چندگانه برآزش دهید.

ب. ضریب تعیین و ضریب تعیین تعدیل شده را به دست آورده و تفسیر کنید.

ج. جدول تحلیل واریانس را به دست آورید و آزمون معنی‌داری مدل را در سطح خطای ۵ درصد انجام دهید.

د. برآورد واریانس مدل را به دست آورده و به کمک آن واریانس پارامترهای مدل را محاسبه کنید.

ه. آزمون  $H_0: \beta_3 = 0$  در برابر  $H_1: \beta_3 \neq 0$  را انجام دهید.

و. بردار مقدارهای برآزش یافته را به دست آورید.

جدول ۱۰.۴: داده‌های درختان کاج

Age	HD	N	MDBH	Age	HD	N	MDBH
۱۹	۵۱/۵	۵۰۰	۷	۱۳	۳۷/۳	۸۰۰	۵/۴
۱۴	۴۱/۳	۹۰۰	۵	۲۱	۵۴/۲	۶۵۰	۶/۴
۱۱	۳۶/۷	۶۵۰	۶/۲	۱۱	۳۲/۵	۵۳۰	۵/۴
۱۳	۳۲/۲	۴۸۰	۵/۲	۱۹	۵۶/۳	۶۸۰	۶/۷
۱۳	۳۹	۵۲۰	۶/۲	۱۷	۵۲/۸	۶۲۰	۶/۷
۱۲	۲۹/۸	۶۱۰	۵/۲	۱۵	۴۷	۹۰۰	۵/۹
۱۸	۵۱/۲	۷۰۰	۶/۲	۱۶	۵۳	۶۲۰	۶/۹
۱۴	۴۶/۸	۷۶۰	۶/۴	۱۶	۵۰/۳	۷۳۰	۶/۹
۲۰	۶۱/۸	۹۳۰	۶/۴	۱۴	۵۰/۵	۶۸۰	۶/۹
۱۷	۵۵/۸	۶۹۰	۶/۴	۲۲	۵۷/۷	۴۸۰	۷/۹

ز. آزمون فرض زیر را در سطح خطای ۵ درصد انجام دهید.

$$H_0: \begin{bmatrix} \beta_1 \\ \beta_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad H_1: \begin{bmatrix} \beta_1 \\ \beta_3 \end{bmatrix} \neq \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

ح. بردار مانده‌ها را با به‌کارگیری رابطه‌ی (۱۷.۴) به‌دست آورید و سپس برآورد واریانس مدل را محاسبه کنید.

۱۵. بر اساس یک مجموعه داده با دو متغیر پیشگوی کمی و متغیر پاسخ متناظر، خلاصه اطلاعات زیر داده شده‌اند.

$$n = 25, \quad y'y = 18310/6290, \quad X'y = \begin{bmatrix} 559/60 \\ 7375/44 \\ 337072 \end{bmatrix},$$

$$(X'X)^{-1} = \begin{bmatrix} 0/1132151 & -0/0044486 & -0/00000836 \\ & 0/0027437 & -0/00000478 \\ & & 0/00000012 \end{bmatrix}$$



مطلوب است

الف. برآورد  $\hat{\beta}$ ب. مقدار آماره‌ی آزمون  $H_0: \beta_2 = 0$ ج. مقدار  $R_{adj}^2$ 

د. برآورد واریانس خطای مدل

۱۶. فرض کنید در یک مدل رگرسیونی با ۴ متغیر مستقل و  $n = 17$  جدول تحلیل واریانس به صورت جدول ۱۱.۴ باشد.

الف. جدول تحلیل واریانس را کامل کنید. ب. ضریب تعیین و  $S_e^2$  را به دست آورید.

جدول ۱۱.۴: جدول تحلیل واریانس تمرین ۱۶

آماره‌ی آزمون	میانگین توان دوم	مجموع توان دوم	درجه‌ی آزادی	منبع تغییرات
۸۵۱/۷۲				مدل
		۲۶۲/۰۷۲		خطا
				کل

۱۷. در یک فرایند پاکسازی ذغال سنگ میزان ذرات معلق بر حسب میلی گرم در لیتر به عنوان متغیر وابسته و مقدار PH فرایند و نوع پلیمر به کار رفته برای انجام واکنش به عنوان متغیرهای مستقل در نظر گرفته شده است. داده‌ها در جدول ۱۲.۴ داده شده است (مایرز، ۱۹۹۰). برای این داده‌ها مدل رگرسیونی را به دست آورده و تأثیر PH بر متوسط ذرات معلق در هر نوع پلیمر را با هم مقایسه کنید.

جدول ۱۲.۴: داده‌های پاک سازی ذغال سنگ

نوع پلیمر	$x_1$ (PH)	$y$ (mg/l)	نوع پلیمر	$x_1$ (PH)	$y$ (mg/l)
p۲	۷/۹	۲۹۷	p۱	۶/۵	۲۹۲
p۲	۸/۷	۳۶۴	p۱	۶/۹	۳۲۹
p۲	۹/۲	۳۷۵	p۱	۷/۸	۳۵۲
p۲	۶/۶	۱۶۷	p۱	۸/۴	۳۷۸
p۲	۷/۰	۲۲۵	p۱	۸/۸	۳۹۲
p۲	۷/۲	۲۴۷	p۱	۹/۲	۴۱۰
p۲	۷/۶	۲۶۸	p۱	۶/۷	۱۹۸
p۲	۸/۷	۲۸۸	p۱	۶/۹	۲۲۷
p۲	۹/۲	۳۴۲	p۱	۷/۵	۲۷۷

۱۸. برای یک مجموعه داده با در نظر گرفتن دو مدل رگرسیونی

$$y_1 = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \varepsilon, \quad y_2 = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

جدول تحلیل واریانس به صورت جدول‌های ۱۳.۴ و ۱۴.۴ به دست آمده است. با به کارگیری این جدول‌ها، آزمون  $H_0: \beta_3 = \beta_4 = 0$  را در سطح خطای ۵ درصد انجام دهید.

جدول ۱۳.۴: جدول تحلیل واریانس مدل اول تمرین ۱۸

p-مقدار	آماره‌ی آزمون	میانگین توان دوم	درجه‌ی آزادی	مجموع توان دوم	منبع پراکندگی
۰/۰۰۰ <sup>a</sup>	۳۰/۲۴۳	۳۵۶۲/۶۱۵	۴	۱۴۲۵۰/۴۶۰	رگرسیون
		۱۱۷/۷۹۸	۱۴۸	۱۷۴۳۴/۱۰۸	خطا
			۱۵۲	۳۱۶۸۴/۵۷	کل

جدول ۱۴.۴: جدول تحلیل واریانس مدل دوم تمرین ۱۸

p-مقدار	آماره‌ی آزمون	میانگین توان دوم	درجه‌ی آزادی	مجموع توان دوم	منبع تغییرات
۰/۰۰۰ <sup>a</sup>	۵۸/۳۹۱	۶۹۳۴/۹۱۳	۲	۱۳۸۶۹/۸۲۵	رگرسیون
		۱۱۸/۷۶۵	۱۵۰	۱۷۸۱۴/۷۴	خطا
			۱۵۲	۳۱۶۸۴/۵۷	کل

۱۹. بر اساس یک مجموعه داده با  $\sum_{i=1}^{32} y_i = ۶۴۲/۹$ ، خلاصه اطلاعات زیر به دست آمده است.

$$y'y = ۱۴۰۴۲/۳۱, \quad x'x = \begin{bmatrix} ۶۴۲۸ \\ ۲۳۸۰/۲۷ \\ ۱۹۰۹/۷۳ \end{bmatrix},$$

$$(x'x)^{-1} = \begin{bmatrix} ۵/۷۶۸۰ & -۱/۱۱۱۱ & -۰/۵۴۰۹ \\ & ۰/۲۲۹۱ & ۰/۰۸۹۲ \\ & & ۰/۰۶۸۴ \end{bmatrix}$$

با توجه به اطلاعات داده شده به پرسش‌های زیر پاسخ دهید.



الف. برآورد بردار ضریب‌های رگرسیونی را به دست آورید.

ب. مقدارهای  $S_{yy}$ ،  $SSR$  و  $SSE$  را به دست آورده و سپس جدول تحلیل واریانس را تشکیل دهید و به کمک آن آزمون معنی‌داری مدل رگرسیونی را در سطح خطای ۵ درصد انجام دهید.

ج. ضریب تعیین مدل را به دست آورید و تفسیر آن را بنویسید.

د. برآورد ماتریس واریانس-کواریانس و برآورد بردار ضریب‌های رگرسیونی را به دست آورید.

ه. برآورد واریانس مقدار برازش یافته‌ی متغیر پاسخ به ازای  $X_t = (1, 3/90, 2/62)$  را به دست آورید.

و. یک بازه‌ی اطمینان ۹۵ درصد با دم‌های برابر برای میانگین متغیر پاسخ به ازای  $X_t$  در قسمت "ه" به دست آورید.

۲۰. در یک مدل رگرسیون خطی چندگانه به صورت  $y = X\beta + \epsilon$  با  $H = X(X'X)^{-1}X'$ ، نشان دهید

$$R^2 = 1 - \frac{y'(I - H)y}{y'y - n\bar{y}^2},$$

که در آن  $n$  تعداد داده‌ها است.

۲۱. در یک مدل رگرسیون خطی با  $p - 1$  متغیر پیشگو و خطاهای  $(\epsilon_1, \dots, \epsilon_n) \stackrel{iid}{\sim} N(0, \sigma^2)$ ، نشان دهید

$$\frac{(\hat{\beta} - \beta)'(X'X)(\hat{\beta} - \beta)}{\sigma^2} \sim \chi_p^2$$

۲۲. اگر در یک مدل رگرسیون خطی با  $p - 1$  متغیر پیشگو، ضریب تعیین را با  $R^2$  و ضریب تعیین مدل کاهش یافته با در نظر گرفتن فرض  $H$  را با  $R_H^2$  نمایش دهیم، نشان دهید آماره‌ی آزمون خطی کلی برای آزمون  $H$  برابر است با

$$F_0 = \frac{(n - p)(R^2 - R_H^2)}{1 - R^2}$$

۲۳. مدل رگرسیون خطی  $y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i$ ،  $i = 1, \dots, n$ ، را در نظر بگیرید، که در آن  $(\epsilon_1, \dots, \epsilon_n) \stackrel{iid}{\sim} N(0, \sigma^2)$  و  $\sigma^2$  نامعلوم باشد. بر اساس داده‌ها، خلاصه اطلاعات زیر به دست

آمده است

$$\sum_{i=1}^{15} x_{i1}^2 = 20, \quad \sum_{i=1}^{15} x_{i1}x_{i2} = 5, \quad \sum_{i=1}^{15} x_{i2}^2 = 10, \quad \sum_{i=1}^{15} x_{i1}y_i = 20,$$

$$\sum_{i=1}^{15} y_i x_{i2} = 10, \quad \sum_{i=1}^{15} y_i^2 = 25$$

الف. برآورد کم‌ترین توان دوم  $\beta_1 + \beta_2$  را بیابید.

ب. برآورد کم‌ترین توان دوم  $\beta_1 + \beta_2$  را با در نظر گرفتن فرض  $\beta_1 = \beta_2 + 0/25$  بیابید.

ج. واریانس  $\hat{\beta}_2$  در قسمت "ب" را بیابید و نشان دهید که از واریانس  $\frac{1}{8}(\hat{\beta}_1 + \hat{\beta}_2) - \frac{1}{8}\bar{y}$  کمتر است.  $\bar{y} = \frac{1}{15}(\hat{\beta}_1 + \hat{\beta}_2) - \frac{1}{8}$  که در آن  $\hat{\beta}_1 + \hat{\beta}_2$  برآوردگر کم‌ترین توان دوم  $\beta_1 + \beta_2$  در قسمت "الف" است، کمتر است.