

رگرسیون خطی ساده - آمار و احتمال مهندسی -

مدرس: مشکانی فراهانی

دانشگاه صنعتی امیرکبیر

۹ دی ۱۳۹۹

از رگرسیون برای به الگو در آوردن رابطه بین متغیرهای آماری استفاده می‌شود.

زیرا هدف بیشتر تحقیق‌ها ارزیابی روابط میان مجموعه‌ای از متغیرهاست.

رگرسیون چیست؟ می‌توان گفت رگرسیون تعیین روابط نادقیق بین متغیرهای آماری و تحلیل این روابط است.

فرض کنید می‌خواهیم بدانیم که آیا مصرف سیگار با متغیرهای اجتماعی مثل سن، تحصیلات، درآمد و قیمت سیگار رابطه دارد یا خیر.

رابطه بین مصرف سیگار با متغیرهای ذکر شده به شکل یک معادله یا الگویی است که یک متغیر وابسته (متغیر پاسخ) را به یک یا چند متغیر مستقل (متغیر پیش‌گو) مربوط می‌کند.

در رگرسیون خطی متغیر پاسخ یک متغیر پیوسته است اما متغیرهای پیش‌گو می‌توانند گسسته یا پیوسته باشند.

متغیر پاسخ را با Y و متغیر پیش‌گو را با X نشان می‌دهند.

در ساده‌ترین روش‌های رگرسیون خطی فرض می‌شود که متغیرها به وسیله یک معادله خط راست در ارتباط هستند.

شکل کلی یک خط راست به صورت $Y = \alpha + \beta X$ است، که در آن α را عرض از مبدأ و β را شیب خط می‌نامند.

الگوی رگرسیون خطی ساده عبارت است از: $Y_i = \alpha + \beta X_i + E_i, \quad i = 1, 2, \dots, n$

در معادله بالا

α عرض از مبدأ و β شیب خط است؛ که ثابت اما مجهول هستند.

α و β را **ضرایب رگرسیونی** می‌نامند.

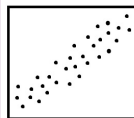
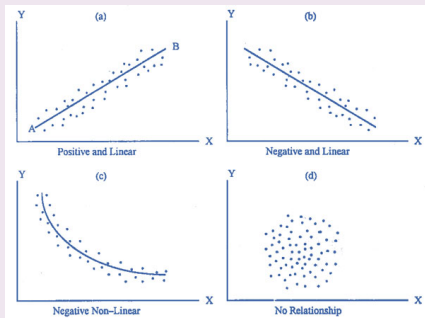
E جمله‌ی خطا است که متغیری تصادفی و غیر قابل مشاهده است: $E(E_i) = 0, \quad Var(E_i) = \sigma^2$

در مدل‌های رگرسیونی رابطه میان متغیر پیشگو و متغیر پاسخ به صورت یک رابطه ریاضی نیست. زیرا به‌ازای هر مقدار از متغیر پیشگو ممکن است چند مقدار برای متغیر پاسخ مشخص شود. بنابراین یک مؤلفه تصادفی در رابطه‌ی آن‌ها وجود دارد E_i .

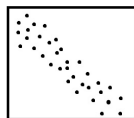
برای به دست آوردن رابطه‌ی بین متغیر پیشگو x و متغیر پاسخ Y ابتدا یک نمونه تصادفی از جامعه‌ی مد نظر جمع‌آوری می‌کنیم.

یعنی به ازای مقادیر x_1, x_2, \dots, x_n از متغیر پیشگو مقادیر متغیر پاسخ $Y|x$ یعنی y_1, y_2, \dots, y_n را اندازه‌گیری می‌کنیم. حال مقادیر مشاهده شده از نمونه‌ی تصادفی یعنی $(x_1, y_1), \dots, (x_n, y_n)$ را داریم.

برای پی بردن به رابطه بین x و $Y|x$ این مشاهدات را به صورت نقاطی در دستگاه مختصات رسم می‌کنیم. این نمودار را **نمودار پراکندگی** یا **پراکنش** می‌نامند.



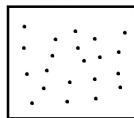
positive linear
association



negative linear
association



nonlinear
association



no association

در صورتی که بین x و y رابطه‌ای وجود داشته باشد، می‌توان یک خط یا منحنی را بر این نقاط عبور داد. این خط یا منحنی رابطه میان x و y را مشخص می‌کند.

این خط یا منحنی را معادله رگرسیونی Y روی x می‌گویند.

منظور از **رگرسیون خطی** آن است که میانگین $Y|x$ به طور خطی با x در ارتباط است؛ یعنی

$$E(Y|x) = \alpha + \beta x$$

در نمونه تصادفی Y_i برابر با $E(Y)$ نیست؛ اختلاف آن‌ها را مقدار خطای E_i نشان می‌دهند.

مقدار مشاهده شده‌ی E_i را با ε_i نشان می‌دهند:

$$y_i = \alpha + \beta x_i + \varepsilon_i$$

برای پیش‌بینی مقادیر $Y|x_i$ از روی خط رگرسیونی باید پارامترهای α و β را با استفاده از مشاهدات برآورد کنیم.

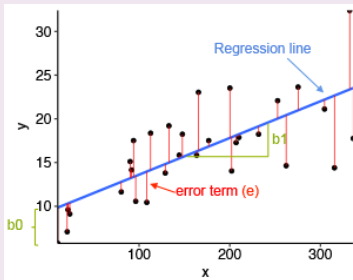
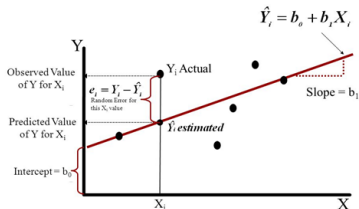
$$\hat{y}_i = \hat{\alpha} + \hat{\beta}x_i \quad \text{مقدار برآورد شده‌ی } \alpha \text{ و } \beta \text{ را با } \hat{\alpha} \text{ و } \hat{\beta} \text{ نشان می‌دهیم:}$$

β و α را ضرایب رگرسیونی می‌نامند.

تفاوت i -امین مقدار مشاهده شده و i -امین مقدار برازش شده‌ی متناظر با آن را باقی‌مانده i -ام می‌نامند:

$$e_i = y_i - \hat{y}_i = y_i - \hat{\alpha} - \hat{\beta}x_i$$

Simple Linear Regression Model



تعیین بهترین خط راست برازش

برای برازش بهترین خط راست از روش کمترین توان‌های دوم خطا استفاده می‌کنیم.

در این روش بهترین خط راست برازش خطی است که مجموع توان دوم خطاها را مینیمم کند.

هرچه انحرافات مقادیر مشاهده شده از این خط کمتر باشد، این خط به داده‌ها نزدیک‌تر است.

پس α و β را طوری برآورد می‌کنیم که $SSE = \sum_{i=1}^n e_i^2$ مینیمم شود.

$$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta}x_i)^2$$

تعیین بهترین خط راست برازش

قضیه: در مدل رگرسیون خطی $Y_i = \alpha + \beta x_i + E_i$ مقادیر $\hat{\alpha}$ و $\hat{\beta}$ که مجموع توان‌های دوم باقی‌مانده‌ها را مینیمم کند، عبارتند از:

$$\hat{\beta} = \frac{s_{xy}}{s_{xx}} \quad \hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

$$s_{xx} = \sum (x_i - \bar{x})^2 = \sum x_i^2 - n\bar{x}^2$$

$$s_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i y_i - n\bar{x}\bar{y}$$

اثبات:

$$\frac{\partial SSE}{\partial \hat{\alpha}} = 0 \Rightarrow -2 \left(\sum y_i - n\hat{\alpha} - \hat{\beta} \sum x_i \right) = 0$$

$$\frac{\partial SSE}{\partial \hat{\beta}} = 0 \Rightarrow -2 \left(\sum x_i y_i - \hat{\alpha} \sum x_i - \hat{\beta} \sum x_i^2 \right) = 0$$

اثبات قضیه

$$\sum y_i - n\hat{\alpha} - \hat{\beta} \sum x_i = 0$$

$$n\hat{\alpha} = \sum y_i - \hat{\beta} \sum x_i$$

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$$

$$\sum x_i y_i - \hat{\alpha} \sum x_i - \hat{\beta} \sum x_i^2 = 0$$

$$\sum x_i y_i - (\bar{y} - \hat{\beta}\bar{x}) \sum x_i - \hat{\beta} \sum x_i^2 = 0$$

$$\hat{\beta} \left(\bar{x} \sum x_i - \sum x_i^2 \right) = \bar{y} \sum x_i - \sum x_i y_i$$

$$\hat{\beta} = \frac{\sum x_i y_i - n\bar{x}\bar{y}}{\sum x_i^2 - n\bar{x}^2} = \frac{s_{xy}}{s_{xx}}$$

مثال ۱

جدول زیر مقادیر فسفر موجود در جیره غذایی ۷ گوساله را همراه با افزایش وزن آن‌ها در طول یک هفته نشان می‌دهد. به این داده‌ها خط رگرسیون برازش دهید.

x	۲۵	۳۰	۳۵	۴۰	۴۵	۵۰	۵۵
y	۴/۱	۴/۳	۴/۳	۵/۵	۵/۹	۶/۵	۶/۸
x_i^2	۶۲۵	۹۰۰	۱۲۲۵	۱۶۰۰	۲۰۲۵	۲۵۰۰	۳۰۲۵
$x_i y_i$	۱۰۲/۵	۱۲۹	۱۵۰/۵	۲۲۰	۲۶۵/۵	۳۲۵	۳۷۴

$$\bar{x} = \frac{\sum x_i}{n} = \frac{۲۵ + \dots + ۵۵}{۷} = ۴۰ \quad \bar{y} = \frac{\sum y_i}{n} = \frac{۴/۱ + \dots + ۶/۸}{۷} = ۵/۳۴$$

$$s_{xx} = \sum x_i^2 - n\bar{x}^2 = ۱۱۹۰۰ - (۷ \times ۴۰^2) = ۷۰۰$$

$$s_{xy} = \sum x_i y_i - n\bar{x}\bar{y} = ۱۵۶۶/۵ - (۷ \times ۴۰ \times ۵/۳۴) = ۷۱/۳$$

$$\hat{\beta} = \frac{s_{xy}}{s_{xx}} = \frac{۷۱/۳}{۷۰۰} = ۰/۱۰۲ \quad \& \quad \hat{\alpha} = \bar{y} - \hat{\beta}\bar{x} = ۵/۳۴ - (۰/۱۰۲ \times ۴۰) = ۱/۲۶$$

$$\Rightarrow \hat{y}_i = ۱/۲۶ + ۰/۱۰۲ x_i \quad \text{معادله خط رگرسیونی:}$$

مثال ۲

یک کمپانی می‌خواهد تأثیر تبلیغات را در فروش کالاهای تولید شده بررسی کند. بدین منظور داده‌های ثبت شده در جدول زیر را بعد از ۱۰ ماه به دست آورد.

الف- معادله خط رگرسیونی را به دست آورید.

ب- نقاط داده شده و خط براورد شده را در یک دستگاه مختصات رسم کنید.

ج- اگر هزینه تبلیغات ۱ باشد، حجم فروش را پیش‌بینی کنید.

ماه	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰
هزینه تبلیغات	۱/۲	۰/۸	۱	۱/۳	۰/۷	۰/۸	۱	۰/۶	۰/۹	۱/۱
حجم فروش	۱۰۱	۹۲	۱۱۰	۱۲۰	۹۰	۸۲	۹۳	۷۵	۹۱	۱۰۰
x_i^2	۱/۴۴	۰/۶۴	۱	۱/۶۹	۰/۴۹	۰/۶۴	۱	۰/۳۶	۰/۸۱	۱/۲۱
$x_i y_i$	۱۲۱/۲	۷۳/۶	۱۱۰	۱۵۶	۶۳	۶۵/۶	۹۳	۴۵	۸۱/۹	۱۱۰

$$\bar{x} = \frac{1/2 + \dots + 1/1}{10} = 0/94$$

$$\bar{y} = \frac{101 + \dots + 100}{10} = 95/4$$

مثال ۲

ادامه راه حل:

$$s_{xx} = \sum x_i^2 - n\bar{x}^2 = 9/28 - (10 \times 0/94^2) = 0/444$$

$$s_{xy} = \sum x_i y_i - n\bar{x}\bar{y} = 924/8 - (10 \times 0/94 \times 95/4) = 28/04$$

$$\hat{\beta} = \frac{s_{xy}}{s_{xx}} = \frac{28/04}{0/444} = 63/15$$

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x} = 95/4 - (63/15 \times 0/94) = 36/04$$

الف- $\hat{y}_i = 36/04 + 63/15 x_i$

ج- $\hat{y}_i = 36/04 + (63/15 \times 1) = 99/19$

به ازای هر واحد افزایش در x ، مقدار برآورد شده y به اندازه $\hat{\beta}$ تغییر می کند.
زمانی که $x = 0$ است، مقدار متوسط y برابر با $\hat{\alpha}$ است.

واریانس $\hat{\beta}$

$$\begin{aligned}
 Var(\hat{\beta}) &= Var\left(\frac{s_{xy}}{s_{xx}}\right) = \frac{1}{s_{xx}^2} Var\left[\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})\right] \\
 &= \frac{1}{s_{xx}^2} \sum_{i=1}^n (x_i - \bar{x})^2 Var(Y_i - \bar{Y}) \\
 &= \frac{1}{s_{xx}^2} \sum_{i=1}^n (x_i - \bar{x})^2 Var(Y_i) \\
 &= \frac{1}{s_{xx}^2} \sum_{i=1}^n (x_i - \bar{x})^2 \sigma^2 \\
 &= \frac{\sigma^2}{s_{xx}}
 \end{aligned}$$

برآورد σ^2

برای برآورد σ^2 از مجموع انحراف‌های Y_i ها حول برآورد آن‌ها در هر سطح از x_i یعنی \hat{Y}_i استفاده می‌کنیم:

$$\hat{\sigma}^2 = S^2 = \frac{SSE}{n-2} = \frac{\sum (Y_i - \hat{Y}_i)^2}{n-2} = \frac{\sum (Y_i - \hat{\alpha} - \hat{\beta}x_i)^2}{n-2} = \frac{s_{yy} - \hat{\beta}s_{xy}}{n-2}$$

اثبات:

$$\begin{aligned} SSE &= \sum_{i=1}^n [Y_i - \hat{\alpha} - \hat{\beta}x_i]^2 = \sum_{i=1}^n [Y_i - \bar{Y} + \hat{\beta}\bar{x} - \hat{\beta}x_i]^2 \\ &= \sum_{i=1}^n [(Y_i - \bar{Y}) - \hat{\beta}(x_i - \bar{x})]^2 \\ &= \sum_{i=1}^n (Y_i - \bar{Y})^2 - 2\hat{\beta} \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y}) + \hat{\beta}^2 \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= s_{yy} - 2\hat{\beta}s_{xy} + \hat{\beta}^2 s_{xx} = s_{yy} - \hat{\beta}s_{xx} \end{aligned}$$

فاصله اطمینان برای $\hat{\beta}$

برای ساختن فاصله اطمینان نیاز به داشتن یک کمیت همراه با تابع توزیع آن داریم.

تاکنون هیچ فرضی روی توزیع متغیرهای تصادفی E_i و Y_i نگذاشتیم.

در مدل رگرسیونی $Y_i = \alpha + \beta x_i + E_i$ چون α ، β و x_i ثابت هستند، با تعیین توزیع E_i می‌توان توزیع Y_i را تعیین کرد.

فرض کنید E_1, \dots, E_n از یکدیگر مستقل بوده و از توزیع نرمال با میانگین صفر و واریانس σ^2 تبعیت کنند:
 $E_i \sim N(0, \sigma^2)$

پس $Y_i = \alpha + \beta x_i + E_i \sim N(\alpha + \beta x_i, \sigma^2)$

از طرفی چون $\hat{\beta} = \frac{s_{xy}}{s_{xx}} = \frac{1}{s_{xx}} \sum (x_i - \bar{x})(Y_i - \bar{Y}) = \frac{1}{s_{xx}} \sum (x_i - \bar{x})Y_i$
ترکیب خطی از متغیرهای تصادفی مستقل نرمال است، پس $\hat{\beta} \sim N(\beta, \frac{\sigma^2}{s_{xx}})$

فاصله اطمینان برای $\hat{\beta}$

با توجه به تعریف توزیع t داریم:

$$\frac{\hat{\beta} - \beta}{S/\sqrt{s_{xx}}} = \frac{\frac{\hat{\beta} - \beta}{\sigma/\sqrt{s_{xx}}}}{\sqrt{\frac{(n-2)S^2}{\sigma^2(n-2)}}} = \frac{N(0, 1)}{\sqrt{\frac{\chi^2_{(n-2)}}{n-2}}} \sim t_{(n-2)}$$

پس یک فاصله اطمینان $(1 - \alpha) \cdot 100\%$ برای β به صورت زیر است:

$$\beta \in \left(\hat{\beta} - t_{1-\frac{\alpha}{2}, (n-2)} \frac{s}{\sqrt{s_{xx}}} , \hat{\beta} + t_{1-\frac{\alpha}{2}, (n-2)} \frac{s}{\sqrt{s_{xx}}} \right)$$

که در آن $\hat{\beta} = \frac{s_{xy}}{s_{xx}}$ و $s^2 = \frac{SSE}{n-2} = \frac{s_{yy} - \hat{\beta}s_{xy}}{n-2}$ برقرار است.

مثال ۳

مواد اولیه‌ای که برای ساختن الیاف مصنوعی به کار می‌روند، در انبار مرطوبی نگهداری می‌شود. نتایج حاصل از اندازه‌گیری رطوبت نسبی در انبار و میزان رطوبت در یک نمونه مواد اولیه (هر دو بر حسب درصد) در ۱۲ روز در جدول زیر آورده شده است.

الف- برآورد خط رگرسیونی را بنویسید.
ب- یک فاصله اطمینان ۹۵ درصدی برای β بسازید.

رطوبت انبار	۴۲	۳۵	۵۰	۴۳	۴۸	۶۲	۳۱	۳۶	۴۴	۳۹	۵۵	۴۸
رطوبت مواد اولیه	۱۲	۸	۱۴	۹	۱۱	۱۶	۷	۹	۱۲	۱۰	۱۳	۱۱

$$s_{xx} = \sum x_i^2 - n\bar{x}^2 = ۸۵۴/۹۱۷$$

$$s_{xy} = \sum x_i y_i - n\bar{x}\bar{y} = ۲۳۰$$

$$\hat{\beta} = \frac{s_{xy}}{s_{xx}} = ۰/۲۶۹$$

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x} = -۰/۹۴۸$$

$$\hat{y} = \hat{\alpha} + \hat{\beta}x = -۰/۹۴۸ + ۰/۲۶۹x$$

مثال ۳

$$s_{yy} = \sum y_i^2 - n\bar{y}^2 = ۷۴$$

$$s^2 = \frac{s_{yy} - \hat{\beta}s_{xy}}{n-2} = ۱/۲۱۳$$

$$1 - \alpha = ۰/۹۵ \Rightarrow 1 - \frac{\alpha}{2} = 1 - \frac{۰/۰۵}{2} = ۰/۹۷۵ \Rightarrow t_{۰/۹۷۵, (12-2)} = ۲/۲۳$$

$$\beta \in \left(\hat{\beta} - t_{1-\frac{\alpha}{2}, (n-2)} \frac{s}{\sqrt{s_{xx}}}, \hat{\beta} + t_{1-\frac{\alpha}{2}, (n-2)} \frac{s}{\sqrt{s_{xx}}} \right)$$

$$\beta \in \left(۰/۲۶۹ - ۲/۲۳ \times \sqrt{\frac{۱/۲۱۳}{۱۵۴/۹۱۷}}, ۰/۲۶۹ + ۲/۲۳ \times \sqrt{\frac{۱/۲۱۳}{۱۵۴/۹۱۷}} \right)$$

$$\beta \in (۰/۱۸۵, ۰/۳۵۳)$$

آزمون فرض برای β

مراحل انجام آزمون فرضیه‌های مربوط به شیب خط رگرسیونی به صورت زیر است:

۱- نوشتن فرضیه‌های آزمون (شبیه به مرحله ۳)

۲- محاسبه آماره آزمون $T_* = \frac{\hat{\beta} - \beta_*}{S/\sqrt{s_{xx}}}$

۳- تعیین ناحیه بحرانی

$$\begin{cases} H_* : \beta = \beta_* \\ H_1 : \beta \neq \beta_* \end{cases}$$

$$\begin{cases} H_* : \beta \leq \beta_* \\ H_1 : \beta > \beta_* \end{cases}$$

$$\begin{cases} H_* : \beta \geq \beta_* \\ H_1 : \beta < \beta_* \end{cases}$$

$$C : |T_*| > t_{1-\frac{\alpha}{2}, (n-2)}$$

$$C : T_* > t_{1-\alpha, (n-2)}$$

$$C : T_* < -t_{1-\alpha, (n-2)}$$

۴- نتیجه‌گیری بر مبنای مرحله ۳

مثال ۴

یک سازنده‌ی مواد افزودنی به بنزین ادعا می‌کند که ماده‌ی افزودنی A در کاهش اکسید نیتروژن خارج شده از خودرو مؤثر است. برای این منظور ۱۰ خودرو از یک مدل تحت آزمایش قرار می‌گیرند. ابتدا میزان اکسید نیتروژن خارج شده از هر ماشین بدون اضافه کردن ماده‌ی A اندازه گرفته می‌شود. سپس به هر یک میزان معینی از ماده A را به باک پر بنزین اضافه کرده و میزان نیتروژن خارج شده از آن خودروها را اندازه‌گیری کرده‌اند. تقلیل در میزان اکسید نیتروژن به عنوان متغیر پاسخ ثبت شده است.

الف- برآورد خط رگرسیونی را به دست آورید.

ب- آیا دلیل قانع‌کننده‌ای برای این وجود دارد که افزایش ماده A باعث کاهش اکسید نیتروژن می‌شود؟
($\alpha = 0.05$)

میزان ماده A	۱	۱	۲	۳	۴	۴	۵	۶	۶	۷
تقلیل در میزان NO	۲/۱	۲/۵	۳/۱	۳	۳/۸	۳/۲	۴/۳	۳/۹	۴/۴	۴/۸

$$s_{xx} = \sum x_i^2 - n\bar{x}^2 = 193 - (10 \times 3/9^2) = 40/9$$

$$s_{yy} = \sum y_i^2 - n\bar{y}^2 = 130/0.5 - (10 \times 3/51^2) = 6/849$$

$$s_{xy} = \sum x_i y_i - n\bar{x}\bar{y} = 152/7 - (10 \times 3/93/51) = 15/81$$

مثال ۴

$$\hat{\beta} = \frac{s_{xy}}{s_{xx}} = ۰/۳۸۷ \quad \hat{\alpha} = \bar{y} - \hat{\beta}\bar{x} = ۲$$

$$\hat{y} = \hat{\alpha} + \hat{\beta}x = ۲ + ۰/۳۸۷ x$$

$$\text{ب-} \quad \begin{cases} H_0 : \beta = ۰ \\ H_1 : \beta > ۰ \end{cases}$$

$$\hat{\sigma}^2 = \frac{SSE}{n-2} = \frac{s_{yy} - \hat{\beta}s_{xy}}{n-2} = \frac{۶/۸۴۹ - (۰/۳۸۷ \times ۱۵/۸۱)}{۱۰-2} = ۰/۰۹۲$$

$$T_0 = \frac{\hat{\beta} - \beta_0}{S/\sqrt{s_{xx}}} = \frac{۰/۳۸۷ - ۰}{\sqrt{\frac{۰/۰۹۲}{۴۰/۹}}} = ۸/۱۴۱$$

$$C : T_0 > t_{1-\alpha, (n-2)} \Rightarrow ۸/۱۴۱ > ۱/۸۶$$

$$\alpha = ۰/۰۵ \Rightarrow t_{۰/۹۵, (۱۰-2)} = ۱/۸۶$$

فرض صفر رد می‌شود؛ یعنی افزایش ماده افزودنی A باعث افزایش تقلیل در میزان اکسید نیتروژن می‌شود.

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.945 ^a	.892	.879	.30365

a. Predictors: (Constant), Material

ANOVA^a

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	6.111	1	6.111	66.284	.000 ^b
	Residual	.738	8	.092		
	Total	6.849	9			

a. Dependent Variable: NitrogenOxides

b. Predictors: (Constant), Material

Coefficients^a

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.
		B	Std. Error	Beta		
1	(Constant)	2.002	.209		9.600	.000
	Material	.387	.047	.945	8.141	.000

a. Dependent Variable: NitrogenOxides

ضریب همبستگی خطی

تا کنون فرض کردیم که متغیر مستقل x یک متغیر کنترل شده است و یک متغیر تصادفی نیست. حال فرض کنید که هم X و هم Y هر دو متغیر تصادفی باشند.

برای سنجش میزان وابستگی دو متغیر تصادفی X و Y از معیاری به نام ضریب همبستگی خطی استفاده می‌شود:

$$\rho = \frac{Cov(X, Y)}{\sqrt{Var(X) Var(Y)}} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

برای برآورد ضریب همبستگی، یک نمونه تصادفی $(X_1, Y_1), \dots, (X_n, Y_n)$ را انتخاب می‌کنیم و از روی این نمونه کوواریانس X و Y ، واریانس X و واریانس Y را به صورت زیر برآورد می‌کنیم:

$$\sigma_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{S_{XX}}{n-1}$$

$$\sigma_Y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 = \frac{S_{YY}}{n-1}$$

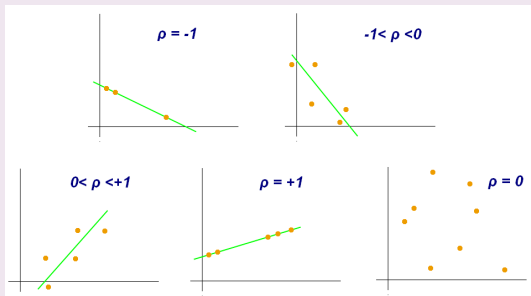
$$\sigma_{XY} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = \frac{S_{XY}}{n-1}$$

ضریب همبستگی خطی

در نتیجه براوردگر ضریب همبستگی خطی به صورت زیر است:

$$\hat{\rho} = R = \frac{S_{XY}}{\sqrt{S_{XX} S_{YY}}}$$

- توجه کنید که همواره $-1 \leq R \leq 1$



مثال ۵

در مثال ۴ ضریب همبستگی نمونه را به دست آورده و آن را تحلیل کنید.

راه حل:

$$\hat{\rho} = R = \frac{s_{xy}}{\sqrt{s_{xx} s_{yy}}} = \frac{۱۵/۸۱}{\sqrt{۴۰/۹ \times ۶/۸۴۹}} = ۰/۹۴۵$$

چون مقدار R به یک نزدیک است پس یک رابطه خطی قوی در جهت مثبت بین X و Y برقرار است.