

فصل ۱۱ تحلیل رگرسیون خطی

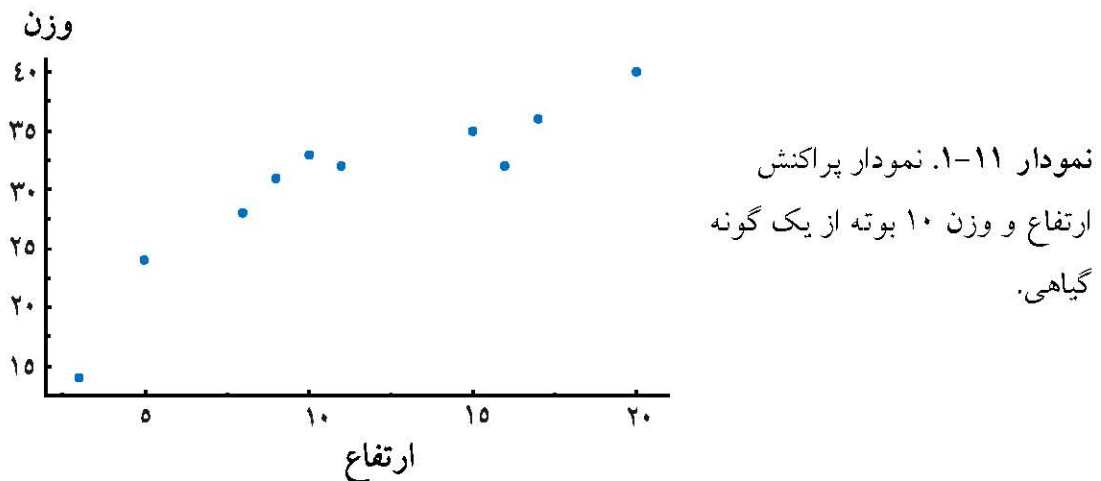
۱-۱۱ مقدمه

یکی از اهداف اغلب تحقیقات علمی بررسی روابط بین متغیرهای مستقل و پاسخ است. در فصل قبل نحوه بررسی اثر یک یا چند عامل بر متغیر پاسخ در قالب تحلیل واریانس، مورد بحث قرار گرفت. در اینگونه مطالعات محقق می‌توانست سطوح مختلف متغیر مستقل را بطور تصادفی به واحدهای آزمایشی نسبت دهد. مثلاً برای بررسی اثر رژیم غذایی بر روی دام‌ها، به هر دام به طور تصادفی یک نوع رژیم غذایی داده می‌شود. این در حالی است که در مورد بعضی متغیرها مانند اثر وزن بیماران بر روی فشار خون، محقق در تعیین وزن بیماران نقشی ندارد. در اینجا متغیر وزن یک متغیر مشاهده‌ای است زیرا مقادیر آن به طور اتفاقی وجود دارد. در این موارد برای توصیف رابطه متغیرهای مشاهده‌ای با متغیر پاسخ از مدل رگرسیونی استفاده می‌شود. مدل رگرسیونی می‌تواند رابطه متغیرهای کمکی و پاسخ را به صورت خطی یا غیر خطی بیان کند. از مدل‌های رگرسیونی برای توصیف داده‌ها، پیش‌بینی و کنترل متغیر پاسخ استفاده می‌شود. قبل از مطالعه تحلیل رگرسیون خطی، نمودار پراکنش و معیارهای کواریانس و ضریب همبستگی که اندازه و جهت رابطه بین دو متغیر را نشان می‌دهند، شرح داده می‌شوند.

نمودار پراکنش

هرگاه مشاهدات برای هر عنصر شامل مقادیر دو متغیر باشد، برای بررسی ارتباط بین دو متغیر می‌توان از نمودار پراکنش استفاده کرد. برای رسم این نمودار یکی از متغیرها روی محور طولی و متغیر دیگر روی محور عرضی در نظر گرفته شده و نقاط با مختصات حاصل از مقادیر دو متغیر رسم می‌شوند. فرض کنید ارتفاع و وزن ۱۰ بوته از یک گونه گیاهی به صورت زوج‌های

(۳, ۱۴)، (۹, ۳۱)، (۱۱, ۳۲)، (۱۰, ۳۳)، (۱۵, ۳۵)، (۲۰, ۴۰)، (۸, ۲۸)، (۵, ۲۴)، (۱۷, ۳۶)، (۱۶, ۳۲) اندازه‌گیری شده است. نمودار پراکنش این داده‌ها در نمودار ۱-۱۱ نشان داده شده است.



۱۱-۲ کواریانس

کواریانس دو متغیر تصادفی x و y به صورت

$$E(x - \mu_x)(y - \mu_y) = E(xy) - \mu_x \mu_y$$

تعریف و با نماد $cov(x, y)$ یا σ_{xy} نشان داده می‌شود. کواریانس معیاری برای سنجش رابطه خطی بین دو متغیر تصادفی است. حاصلضرب $(x - \mu_x)(y - \mu_y)$ ممکن است مثبت یا منفی باشد، پس امید ریاضی آن یعنی σ_{xy} بر حسب اینکه مثبت یا منفی باشد، بیان می‌کند که x و y به طور متوسط و بصورت خطی در یک جهت یا در دو جهت مخالف حرکت می‌کنند. کواریانس در حالت کلی پارامتری مربوط به جامعه است و تنها با در دست داشتن توزیع احتمالی توأم متغیرهای تصادفی (x, y) قابل محاسبه است. از این رو در کاربردها از برآورد نمونه‌ای آن استفاده می‌شود. اگر نمونه‌ای تصادفی از جامعه به صورت $(x_1, y_1), \dots, (x_n, y_n)$ در اختیار باشد، در اینصورت

$$s_{xy} = \hat{\sigma}_{xy} = \frac{SP_{xy}}{n-1} = \frac{\sum (x - \bar{x})(y - \bar{y})}{n-1}$$

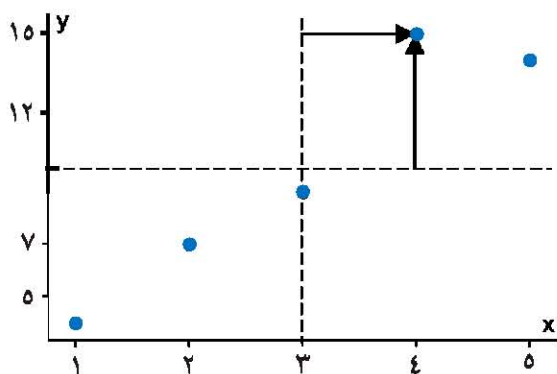
برآوردگری نأریب برای σ_{xy} است. فرض کنید در آزمایشی مقادیر محصول به ازای مقادیر مختلف کود به صورت جدول ۱-۱۱ بوده است.

جدول ۱-۱۱. مقادیر مشاهده شده عملکرد تحت مقادیر مختلف کود.

x (میزان کود)	۱	۲	۳	۴	۵
y (میزان عملکرد)	۴	۷	۹	۱۵	۱۴

کواریانس بین دو متغیر x (میزان کود) و y (میزان محصول) برای نمونه مورد بررسی به صورت زیر برآورد می‌شود.

$$s_{xy} = \frac{\sum (x - \bar{x})(y - \bar{y})}{n - 1} = \frac{\sum xy - \frac{(\sum x)(\sum y)}{n}}{n - 1} = \frac{175 - \frac{15 \times 49}{5}}{4} = 7$$



نمودار ۱۱-۲. نمودار پراکنش مقادیر مشاهده شده محصول (y) تحت مقادیر مختلف کود (x).

پراکنش نقاط مربوطه در این مثال نشان می‌دهد که دو متغیر گرایش به همراهی با هم دارند (نمودار ۱۱-۲). در این مثال حاصلضرب $(x - \mu_x)(y - \mu_y)$ برای هر کدام از نقاط مقداری مثبت به دست می‌آید. در نتیجه اگر مقدار $(x - \mu_x)(y - \mu_y)$ برای همه نقاط محاسبه و با هم جمع شوند، میانگین آنها که کواریانس نام دارد، می‌تواند به عنوان ملاکی برای میزان هم‌تغییری خطی x و y مد نظر باشد. مثبت بودن کواریانس در این مثال نشان می‌دهد که متغیرهای x و y بصورت خطی در یک راستا حرکت می‌کنند.

کواریانس شاخصی است که ارتباط خطی بین دو متغیر را منعکس می‌کند و بصورت

$$\begin{aligned}\sigma_{xy} &= E(x - \mu_x)(y - \mu_y) \\ &= E(xy) - \mu_x\mu_y\end{aligned}$$

تعریف می‌شود. برآورد نمونه‌ای کواریانس بصورت زیر است.

$$\hat{\sigma}_{xy} = s_{xy} = \frac{SP_{xy}}{n-1} = \frac{\sum (x - \bar{x})(y - \bar{y})}{n-1}$$

در عمل به جای کواریانس از ضریب همبستگی پیرسون برای نشان دادن میزان ارتباط بین دو متغیر استفاده می‌شود. زیرا کواریانس به واحد اندازه‌گیری وابسته بوده و با توجه به میانگین نمونه‌ها می‌تواند هر مقداری را اختیار کند.

۱۱-۳ ضریب همبستگی پیرسون

ضریب همبستگی پیرسون یا ضریب همبستگی خطی شدت و جهت همبستگی خطی بین دو متغیر تصادفی را نشان می‌دهد و با اعمال تغییراتی در تعریف کواریانس بصورت

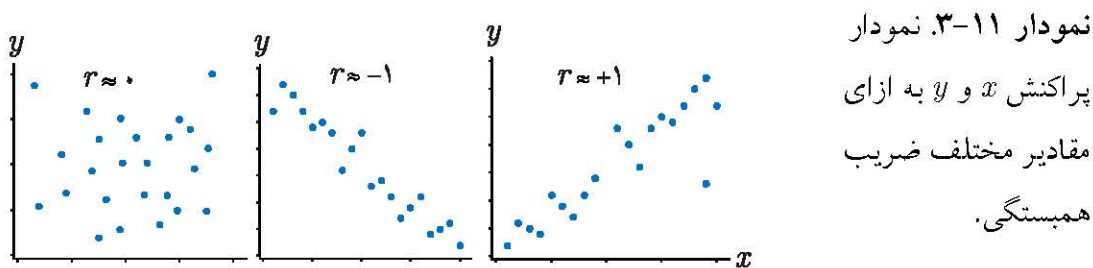
$$\rho(x,y) = \frac{\sigma_{xy}}{\sigma_x\sigma_y}$$

تعریف می‌شود. ضریب همبستگی برخلاف کواریانس به واحد اندازه‌گیری وابسته نیست و مقادیر آن محدود به فاصله (۱ و -۱) است. ضریب همبستگی نیز کمیتی مربوط به جامعه است، از اینرو در کاربردها از برآورد نمونه‌ای آن استفاده می‌شود. اگر $(x_1, y_1), \dots, (x_n, y_n)$ نمونه‌ای تصادفی از جامعه باشد، در اینصورت

$$\eta(x,y) = \hat{\rho}(x,y) = \frac{s_{xy}}{s_x s_y}$$

برآوردگری مناسب برای $\rho(x,y)$ است. مقدار ضریب همبستگی نزدیک به +۱ نشان‌دهنده وجود ارتباط خطی مثبت بین x و y است، یعنی با افزایش x مقدار y نیز افزایش می‌یابد و برعکس. در صورتی که r نزدیک به -۱ باشد، بین دو متغیر ارتباط خطی منفی وجود داشته و با افزایش یکی

از آنها دیگری کاهش می‌یابد. مقادیر r نزدیک به صفر نیز بیانگر عدم وجود ارتباط خطی بین دو متغیر است (نمودار ۱۱-۳).



مثال ۱۱-۱. در آزمایشی به منظور بررسی ارتباط جوانه‌زنی بذر با میزان دما، درصد جوانه‌زنی در اثر ۵ دمای مختلف پس از مدت مشخصی به دست آمده است. فرض کنید نتایج نمونه مورد بررسی به صورت زیر بوده است.

دما (x)	۱	۲	۳	۴	۵
درصد جوانه‌زنی (y)	۱	۱	۲	۲	۴

مقدار r برای داده‌های مثال ۱۱-۱ به صورت زیر به دست می‌آید.

$$r = \frac{SP_{xy}}{\sqrt{SS_x SS_y}}$$

که در آن:

$$SP_{xy} = \sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n} = 37 - \frac{10 \times 10}{5} = 37 - 20 = 7$$

$$SS_x = \sum x_i^2 - \frac{(\sum x_i)^2}{n} = 55 - \frac{10^2}{5} = 55 - 20 = 35$$

$$SS_y = \sum y_i^2 - \frac{(\sum y_i)^2}{n} = 26 - \frac{10^2}{5} = 26 - 20 = 6$$

لذا r برابر با

$$r = \frac{7}{\sqrt{35 \times 6}} = 0.404$$

خواهد بود. این مقدار که از روی نمونه محاسبه شده است تخمینی از ضریب همبستگی جامعه (ρ) می‌باشد. ممکن است آزمون فرض $H_0: \rho = 0$ در مقابل $H_1: \rho \neq 0$ مد نظر باشد. بر اساس فرض صفر، متغیرهای x و y بصورت خطی ناهمبسته‌اند، یعنی بین آنها رابطه خطی وجود ندارد. البته به این معنی نیست که بین این دو متغیر هیچگونه رابطه دیگری نیز وجود ندارد. در صورتی که مقدار ρ نزدیک به $+1$ یا -1 باشد، توزیع آماره r متقارن نبوده و از این رو نمی‌توان با استفاده از الگوهایی نظیر توزیع t و z اقدام به آزمون فرض نمود. فقط در حالتی که ρ صفر باشد، توزیع آماره r ، یک توزیع نرمال با میانگین صفر و واریانس $(1-r^2)/(n-2)$ خواهد بود. پس اگر آزمون $H_0: \rho = 0$ مد نظر باشد، آماره t به صورت

$$t = \frac{r - \rho}{\sigma_r} = \frac{r}{\sqrt{\frac{1-r^2}{n-2}}}$$

محاسبه شده و با t جدول با درجه آزادی $n-2$ مقایسه می‌شود. اگر ρ برابر با صفر نباشد، توزیع r نامتقارن و دارای کجی خواهد بود. پس در مواردی که فرضی غیر از $H_0: \rho = 0$ مورد آزمون باشد، نمی‌توان از آماره t یا z در آزمون فرض استفاده کرد. در چنین مواردی ابتدا با استفاده از تبدیل $z' = 0.5 \ln \frac{1+r}{1-r}$ به متغیر جدیدی با توزیع نرمال تبدیل می‌شود تا امکان انجام آزمون با استفاده از توزیع z به وجود آید. سپس آزمون‌های آماری بر روی متغیر نرمال حاصل z' ، انجام می‌شود. ثابت می‌شود که z' دارای توزیع تقریباً نرمال با انحراف معیار $\sqrt{1/(n-3)}$ می‌باشد. همچنین می‌توان ابتدا فاصله اطمینان $(1-\alpha)$ را برای z' محاسبه نموده و سپس با استفاده از رابطه زیر آن را به فاصله اطمینان r برگرداند.

$$r = \frac{e^{2z'} - 1}{e^{2z'} + 1}$$

مثال ۱۱-۲. ضریب همبستگی بین درصد پروتئین و درصد چربی، در ۵۰ مورد اندازه‌گیری روی بذر سویا، برابر با ۰/۵۲۷ برآورد شده است. فرض $H_0: \rho = 0/5$ را در مقابل $H_1: \rho \neq 0/5$ در سطح ۰/۰۵ بیازمایید.

با استفاده از رابطه $z' = 0/5 \ln \frac{1+r}{1-r}$ مقادیر تبدیلی بصورت $z'_{0/527}$ و $z'_{0/5}$ بدست می‌آیند، لذا

$$z = \frac{z'_r - z'_\rho}{\sigma_{z'}} = \frac{z'_{0/527} - z'_{0/5}}{\sqrt{\frac{1}{n-3}}} = \frac{0/586 - 0/549}{\sqrt{\frac{1}{47}}} = 0/2$$

چون مقدار آماره z از z جدول (۱/۹۶) کوچکتر است، دلیل کافی برای رد فرض $\rho = 0/5$ در سطح ۰/۰۵ وجود ندارد.

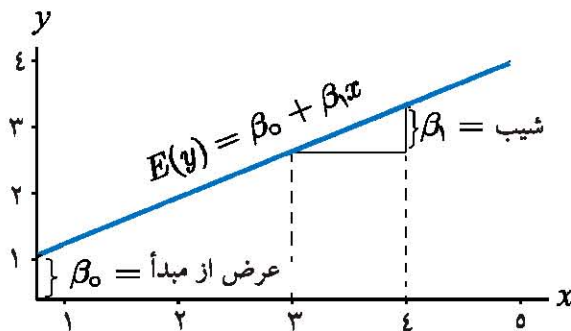
همچنین می‌توان فاصله اطمینان $1 - \alpha$ مثلاً ۰/۹۵ را برای z' بصورت

$$z'_r \pm z_{\alpha/2} \sigma_{z'} = 0/586 \pm 1/96(0/146) = 0/586 \pm 0/286$$

محاسبه کرد. بدین معنی که با احتمال ۰/۹۵، z' در فاصله ۰/۳ تا ۰/۸۸ قرار دارد. با تبدیل z' به r می‌توان حدود اعتماد ۰/۹۵ را برای r به دست آورد. مقدار r برای $z' = 0/3$ ، برابر با ۰/۲۹ و برای $z' = 0/88$ برابر با ۰/۷ می‌باشد. بنابراین با احتمال ۰/۹۵، فاصله ۰/۲۹ تا ۰/۷۰ ضریب همبستگی جامعه را در بر دارد.

۱۱-۴ رگرسیون خطی ساده

رگرسیون خطی ساده ارتباط خطی دو متغیر را بررسی می‌کند. در رگرسیون خطی ساده رابطه بین متغیر مستقل x و متغیر پاسخ y توسط معادله یک خط مستقیم در محور مختصات نشان داده می‌شود. معادله چنین خطی به صورت $E(y) = \beta_0 + \beta_1 x$ نوشته می‌شود که در آن، β_0 عرض از مبدأ خط بوده و نقطه‌ای است که در آن، خط رگرسیون محور y را قطع می‌کند. β_1 نیز شیب خط رگرسیون است و بیانگر میزان افزایش در متغیر y به ازای یک واحد افزایش در متغیر x می‌باشد.



نمودار ۱۱-۴. مدل رگرسیون خطی ساده.

به عنوان مثال درصد جوانه‌زنی بذر با افزایش دما رابطه مستقیم دارد. چون میزان جوانه‌زنی (y) به دما بستگی دارد، آن را متغیر پاسخ یا وابسته می‌نامند. میزان جوانه‌زنی به ازای یک سطح دما در تکرارهای مختلف ثابت نبوده و به دلیل تأثیر عوامل ناشناخته، مقداری متغیر است. اگر میزان دما به طور دلخواه توسط محقق تعیین شود، از نوع داده‌های آزمایشگاهی و در صورتی که محقق نقشی در تعیین آن نداشته باشد، از نوع داده‌های مشاهده‌ای خواهد بود که در هر دو صورت متغیر دما را متغیر مستقل می‌نامند (برای مثال داده‌های مربوط به دما و میزان جوانه‌زنی برای چند گیاه به طور اتفاقی ثبت شده یا میزان جوانه‌زنی یک گونه گیاهی در مناطق با دماهای مختلف اندازه‌گیری شده و محقق می‌خواهد اثر دما را بر آن بررسی کند).

مثال ۱۱-۳. فرض کنید در یک آزمایش، درصد جوانه‌زنی در اثر ۵ دمای مختلف پس از مدت مشخصی به صورت زیر بوده است.

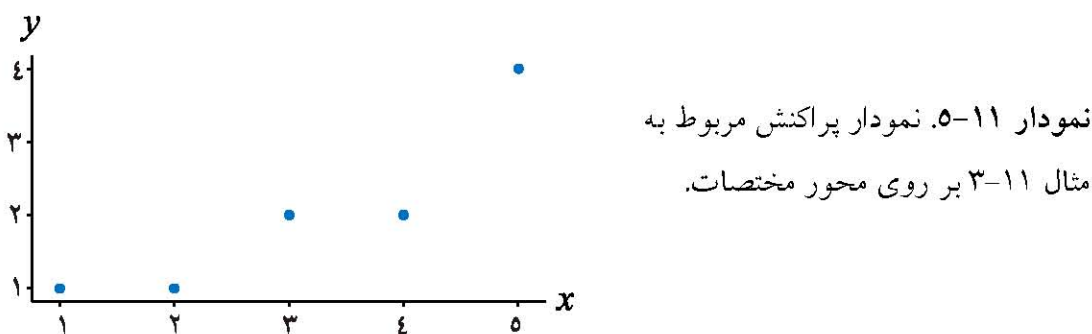
	۱	۲	۳	۴	۵
دما (x)	۱	۲	۳	۴	۵
درصد جوانه‌زنی (y)	۱	۱	۲	۲	۴

نمودار پراکنش درصد جوانه‌زنی به ازای دماهای مختلف (نمودار ۱۱-۵) نشان می‌دهد که دما بر روی میزان جوانه‌زنی اثر دارد. چگونگی این اثر به وسیله معادله‌ای که x و y را به هم ربط می‌دهد توصیف می‌شود. برآورد یک معادله، هم‌ارز برآوردن یک خط به این نقاط است، که خط رگرسیون y بر حسب x نامیده می‌شود و می‌تواند در پیش‌بینی درصد جوانه‌زنی (y) به ازای دماهای مختلف (x) بکار رود. معادله کلی خط رگرسیونی برای میانگین درصد جوانه‌زنی در جامعه به صورت $E(y) = \beta_0 + \beta_1 x$ نشان داده می‌شود که در آن β_0 و β_1 مقادیری ثابت هستند. بدیهی است که مقدار مشاهده شده متغیر وابسته (درصد جوانه‌زنی) به ازای یک مقدار

مشخص دما (متغیر مستقل) دقیقاً برابر با مقدار پیش‌بینی شده نبوده و مقداری خطا وجود خواهد داشت. لذا مدل رگرسیون را می‌توان به صورت

$$y = \beta_0 + \beta_1 x + \varepsilon$$

نیز نوشت، که در آن ε نشان‌دهنده خطا در تعیین y بر اساس x است. پارامترهای β_0 و β_1 ضرایب رگرسیونی نامیده شده و بر اساس مشاهدات نمونه برآورد می‌شوند.



۱۱-۵ برازش خط رگرسیونی

برازش یک مدل رگرسیونی به مجموعه‌ای از مشاهدات، به معنی برآورد پارامترهای مدل بر اساس آن مجموعه از مشاهدات است. خطوط مختلفی را می‌توان از بین نقاط مربوط به نمونه مورد بررسی در مثال ۱۱-۳ عبور داد، ولی مناسب‌ترین آنها خطی است که به ازای آن، مجموع توان‌های دوم خطا حداقل باشد (نمودار ۱۱-۶). معادله خط مذکور که بر اساس اطلاعات نمونه به دست می‌آید، به صورت $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ نشان داده می‌شود. \hat{y} برآورد مقدار y به ازای هر x و $\hat{\beta}_1$ برآورد β_1 است. اگر $(x_1, y_1), \dots, (x_n, y_n)$ نمونه‌ای تصادفی باشد، برآورد مجموع توان‌های دوم خطا بصورت

$$SS_E = \sum (y_i - \hat{y})^2 = \sum [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x)]^2$$

و برآورد ضریب رگرسیون بصورت

$$\hat{\beta}_1 = \frac{SP_{xy}}{SS_x}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

می‌باشد. در مثال ۱۱-۳، این مقادیر بصورت

$$SP_{xy} = \sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n} = 37 - \frac{10 \times 15}{5} = 37 - 30 = 7$$

$$SS_x = \sum x_i^2 - \frac{(\sum x_i)^2}{n} = 55 - \frac{10^2}{5} = 55 - 40 = 15$$

$$\hat{\beta}_1 = \frac{SP_{xy}}{SS_x} = \frac{7}{15} = 0.467$$

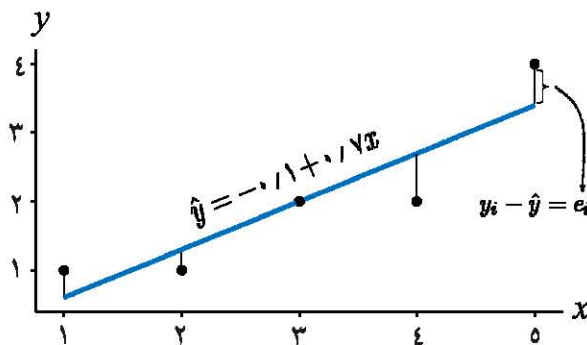
$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = \frac{\sum y_i}{5} - \hat{\beta}_1 \frac{\sum x_i}{5} = \frac{15}{5} - 0.467 \left(\frac{10}{5}\right) = 3 - 0.467(2) = 2.066$$

بدست می‌آیند. لذا معادله خط رگرسیون برآورد شده به صورت

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x = 2.066 + 0.467x$$

است. با استفاده از این معادله می‌توان مقدار y را به ازای هر مقدار مشخص x پیش‌بینی کرد (نمودار ۱۱-۶). مثلاً به ازای $x = 2$ مقدار y به صورت زیر پیش‌بینی خواهد شد.

$$\hat{y} = 2.066 + 0.467x = 2.066 + 0.467(2) = 2.934$$



نمودار ۱۱-۶. خط برازنده شده به داده‌های مثال ۱۱-۳.

۱۱-۶ فرض‌های اولیه تحلیل رگرسیونی

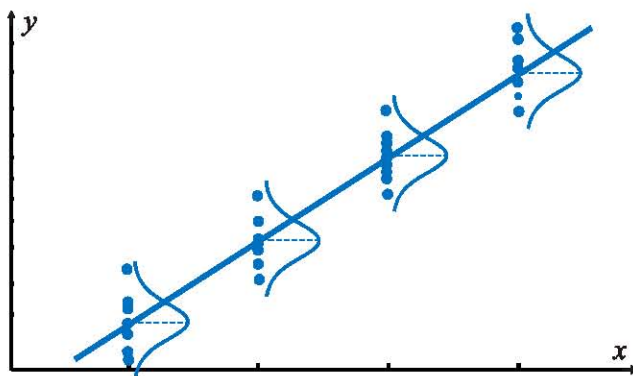
در تحلیل رگرسیونی، برآورد پارامترهای مدل رگرسیونی بر اساس روش حداقل مجموع توان‌های دوم با در نظر گرفتن فرضیات اولیه‌ای در مورد جمله خطا صورت می‌پذیرد. این فرضیات را می‌توان به صورت زیر بیان کرد.

۱- میانگین توزیع خطاها برابر با صفر است. به این معنی که خطاهای مثبت و منفی همدیگر را خنثی کرده و در صورت تکرار آزمایش به ازای یک مقدار مشخص x ، میانگین مقادیر y برابر با $E(y) = \beta_0 + \beta_1 x$ خواهد بود.

۲- واریانس خطا به ازای هر مقدار x ثابت و برابر با مقداری مانند σ^2 است.

۳- خطاها از یکدیگر و از x ناهمبسته می‌باشند. به عبارت دیگر مقدار خطا (ε) به ازای یک مقدار y هیچ اثری بر روی خطای مرتبط با مقدار دیگری از y ندارد (نمودار ۷-۱۱).

بطور خلاصه، خطاها متغیرهایی تصادفی و ناهمبسته با میانگین صفر و واریانس ثابت تلقی می‌شوند. نکته قابل توجه آن است که برای برآورد پارامترها، فرض نرمال بودن خطاها ضروری نیست. زیرا روش کمترین توان‌های دوم روشی جبری است و به توزیع خاصی وابسته نیست. البته برای استنباط در مورد ضرایب رگرسیونی یعنی برای انجام آزمون فرض و ساختن فاصله اطمینان، فرض نرمال بودن خطاها ضرورت پیدا می‌کند.



نمودار ۷-۱۱. توزیع احتمال خطا به ازای مقادیر مختلف متغیر مستقل.

۷-۱۱ برآورد خطای مدل

هر اندازه تغییرات داده‌ها نسبت به خط رگرسیون برآورد شده بیشتر باشد، خطاهای تصادفی و در نتیجه واریانس خطای مدل رگرسیونی بیشتر شده و برآوردهای β_0 و β_1 و همچنین مقدار پیش‌بینی شده \hat{y} به ازای یک مقدار مشخص x دقت کمتری خواهند داشت. به همین دلیل است که واریانس خطا (σ^2) در رابطه برآورد فاصله اطمینان و نیز محاسبه آماره آزمون به چشم می‌خورد. غالباً در عمل σ^2 نامعلوم است و لذا از طریق نمونه برآورد می‌شود. این برآورد که با s^2

نشان داده می‌شود، برابر است با مجموع توان دوم فواصل نقاط از خط رگرسیون تقسیم بر درجه آزادی. از آنجا که در ایجاد مدل رگرسیونی، دو پارامتر (β_0 و β_1) باید برآورد شود، $n - 2$ درجه آزادی برای خطا باقی می‌ماند، لذا:

$$s^2 = \frac{SS_E}{df_E} = \frac{SS_E}{n-2} = \frac{\sum (y_i - \hat{y}_i)^2}{n-2}$$

$$SS_E = \sum (y_i - \hat{y}_i)^2 = SS_y - \hat{\beta}_1 SP_{xy} = SS_y - \hat{\beta}_1^2 SS_x$$

$$SS_y = \sum y_i^2 - \frac{(\sum y_i)^2}{n}$$

در مثال ۱۱-۳ در مورد میزان دما (x) و میزان جوانه‌زنی (y)، مجموع توان‌های دوم خطا به صورت زیر به دست می‌آید.

x	y	$\hat{y} = -0.1 + 0.7x$	$y_i - \hat{y}$	$(y_i - \hat{y})^2$
۱	۱	۰/۶	۴/۰	۰/۱۶
۲	۱	۱/۳	-۰/۳	۰/۰۹
۳	۲	۲/۰	۰	۰
۴	۲	۲/۷	-۰/۷	۰/۴۹
۵	۴	۳/۴	۰/۶	۰/۳۶
مجموع خطاها = ۰				$SS_E = 1/1$

$$s^2 = \frac{SS_E}{n-2} = \frac{1/1}{3} = 0.3333$$

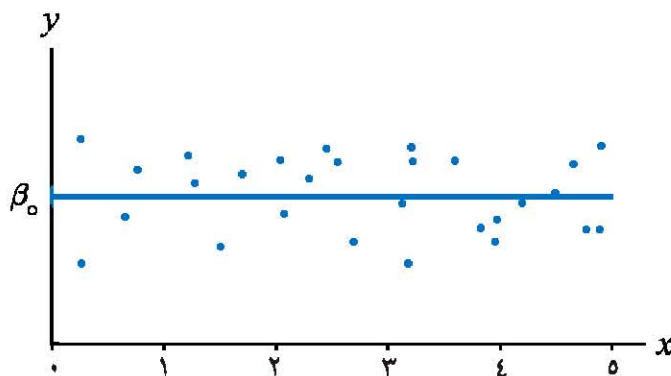
انحراف معیار خطا که خطای معیار مدل رگرسیونی نیز نامیده می‌شود، برابر است با $s = \sqrt{0.3333} = 0.577$ و نشان می‌دهد که در صورت نرمال بودن توزیع خطاها، حدود ۹۵٪ مقادیر مشاهده شده y در فاصله $2s$ از مقادیر پیش‌بینی شده قرار می‌گیرند. در مثال ۱۱-۳ همه مقادیر y در فاصله (2×0.577) از \hat{y} قرار گرفته‌اند.

جدول ۱۱-۲. حل مثال ۱۱-۳ و خروجی مربوطه در برنامه Minitab.

مراحل:					
۱. مشاهدات مربوط به متغیر مستقل (x) و وابسته (y) را در ستون‌های اول و دوم صفحه داده‌ها وارد کنید.					
۲. وارد مسیر مقابل شوید: Stat > Regression > Regression...					
۳. در مقابل گزینه Response ستون دوم و در مقابل گزینه Predictors ستون اول (متغیر مستقل) را انتخاب کنید.					
۴. بر روی OK کلیک کنید.					
The regression equation is $y = -0.100 + 0.700 x$					
Predictor	Coef	SE Coef	T	P	
Constant	-0.1000	0.6351	-0.16	0.885	
x	0.7000	0.1915	3.66	0.035	
S = 0.605530 R-Sq = 81.7% R-Sq(adj) = 75.6%					
Analysis of Variance					
Source	DF	SS	MS	F	P
Regression	1	4.9000	4.9000	13.36	0.035
Residual Error	3	1.1000	0.3667		
Total	4	6.0000			

۱۱-۸ بررسی کارایی مدل رگرسیونی

در عمل ضرایب رگرسیونی تنها بر اساس یک نمونه تصادفی برآورد می‌شوند. لذا اعتماد به این ضرایب بدون در نظر گرفتن توزیع برآوردگرهای مربوط به آنها صحیح نیست. به عبارتی تنها پس از انجام آزمون فرض و کسب اطمینان از معنی‌داری این ضرایب، می‌توان مدل بدست آمده را به جامعه تحت بررسی تعمیم داد. به کمک آزمون فرض‌ها می‌توان مؤثر بودن یا نبودن مدل رگرسیونی در پیش‌بینی مقادیر متغیر وابسته بر اساس مقادیر متغیر مستقل را آزمود. اگر متغیر x در مدل $y = \beta_0 + \beta_1 x + \varepsilon$ نقشی در پیش‌بینی مقادیر y نداشته باشد، شیب خط رگرسیون برابر با صفر خواهد بود، زیرا در این صورت میانگین y در اثر افزایش یا کاهش x تغییری نخواهد کرد (نمودار ۱۱-۸).



نمودار ۱۱-۸ مدل رگرسیونی انجام آزمون بدون اطلاع از توزیع برآوردگرهای ضرایب رگرسیونی (β_1 و β_0) ممکن نیست، در این مرحله، فرض نرمال بودن مانده‌ها نیز به فرضیات اولیه مدل رگرسیونی افزوده می‌شود. در این صورت $\hat{\beta}_1$ دارای توزیع نمونه‌ای نرمال با میانگین β_1 و انحراف معیار $\sigma_{\hat{\beta}_1} = \sigma / \sqrt{SS_x}$ خواهد بود. چون $\sigma_{\hat{\beta}_1}$ از روی نمونه برآورد می‌شود، لذا فرض‌های مذکور را می‌توان به کمک آماره

$$t = \frac{\hat{\beta}_1 - \beta_1}{s_{\hat{\beta}_1}} = \frac{\hat{\beta}_1 - \beta_1}{\frac{s}{\sqrt{SS_x}}}$$

مورد آزمون قرار داد، که دارای توزیع t با $n-2$ درجه آزادی است. در صورتی که قدر مطلق آماره t از t جدول برای سطح معنی‌داری مورد نظر و درجه آزادی $n-2$ بزرگتر باشد فرض صفر رد خواهد شد.

مثال ۱۱-۴. در مثال ۱۱-۳، فرض $H_0: \beta_1 = 0$ را در مقابل $H_1: \beta_1 \neq 0$ در سطح $\alpha = 0.05$ بیازمایید.

$$t = \frac{0.7 - 0}{\frac{0.6055}{\sqrt{10}}} = \frac{0.7}{0.19} = 3.65$$

چون قدر مطلق مقدار آماره t از $t_{0.025, 3} = 3.182$ بزرگتر است، پس شیب خط رگرسیون برابر با صفر نیست. فاصله اطمینان ۹۵٪ برای شیب خط به صورت

$$\hat{\beta}_1 \pm t_{\alpha/2} s_{\hat{\beta}_1} = 0.7 \pm 3.182(0.19) = 0.7 \pm 0.6$$

است. بدین معنی که فاصله ۰/۱ تا ۱/۳ با احتمال ۰/۹۵ شیب واقعی خط رگرسیون جامعه یعنی (β_1) را در بر دارد.

۹-۱۱ ضریب تعیین

همواره قسمتی از تغییرات کل، $SS_y = \sum (y_i - \bar{y})^2$ ، مربوط به خطا یعنی $SS_E = \sum (y_i - \hat{y})^2$ بوده و در صورتی که همه نقاط بر روی خط رگرسیون برآورد شده قرار گیرند، مجموع توان‌های دوم خطا صفر خواهد شد. لذا SS_y را می‌توان به دو بخش SS_E و $SS_y - SS_E$ تقسیم نمود. مقدار $SS_y - SS_E$ در واقع بخشی از تغییرات کل است که توسط مدل رگرسیون تبیین می‌شود و به همین دلیل مجموع توان‌های دوم رگرسیون (SS_R) نام دارد. ضریب تعیین (R^2) بیان می‌کند که چه مقدار از کل تغییرات توسط مدل رگرسیونی تبیین می‌شود. هر چه این مقدار بیشتر باشد رابطه بین x و y قوی‌تر است. ضریب تعیین در رگرسیون خطی ساده برابر است با مربع ضریب همبستگی و بصورت

$$R^2 = \frac{SS_y - SS_E}{SS_y} = \frac{SS_R}{SS_y}$$

برآورد می‌شود. مقدار ضریب تعیین در مثال ۱۱-۳ به صورت

$$R^2 = \frac{SS_y - SS_E}{SS_y} = \frac{6 - 11}{6} = 0.816$$

بدست می‌آید. پس می‌توان گفت که ۰/۸۱ تغییرات جوانه‌زنی (متغیر y) توسط دما (متغیر x) توجیه می‌شود.

۱۱-۱۰ بررسی کارایی مدل رگرسیونی با آزمون F

همانگونه که در نمودار ۱۱-۵ نشان داده شده است، اگر فرض $H_0: \beta_1 = 0$ درست باشد خط رگرسیون افقی بوده و در نتیجه مجموع توان‌های دوم رگرسیون برابر با صفر خواهد بود. به عبارت دیگر میانگین توان‌های دوم رگرسیون نیز برآوردی از واریانس خطا (σ^2) می‌باشد. اگر H_0 درست نباشد مجموع توان‌های دوم رگرسیون برآوردی از واریانس خطا به اضافه واریانس رگرسیون خواهد بود. برای آزمون فرض $H_0: \beta_1 = 0$ با استفاده از توزیع F کافی است که نسبت واریانس رگرسیونی به واریانس خطا محاسبه شده و با F جدول مربوط به درجات آزادی ۱ (مربوط به رگرسیون) و $n - 2$ (مربوط به خطا) مقایسه شود. تأیید بزرگتر بودن واریانس رگرسیون نسبت به واریانس خطا بیانگر معنی‌دار بودن رگرسیون یا ضریب رگرسیون خواهد بود.

در مثال ۱۱-۳ میانگین توان‌های دوم خطا برابر با 0.3667 بدست آمد که از تقسیم توان‌های دوم خطا ($1/1$) بر درجه آزادی آن ($n-2=3$) بدست آمده است. چون رابطه $S_{yy} = SS_R + SS_E$ برقرار است، بنابراین مجموع توان‌های دوم رگرسیون در مثال ۱۱-۳ برابر با $4/9 = 1/1 - 6$ است. البته SS_R را می‌توان مستقیماً بصورت

$$\begin{aligned} \sum (\hat{y} - \bar{y})^2 &= \hat{\beta}_1^2 \sum (x_i - \bar{x})^2 = \hat{\beta}_1^2 SS_x = \hat{\beta}_1^2 SP_{xy} \\ &= (0.7)^2 \times 10 = 0.7 \times 7 = 4/9 \end{aligned}$$

محاسبه نمود. با توجه به اینکه SS_R رگرسیون دارای یک درجه آزادی است، SS_R برابر با MSR خواهد بود. لذا نسبت دو واریانس تحت عنوان آماره F برابر با

$$F = \frac{MSR}{MSE} = \frac{4/9}{0.3667} = 13/36$$

است. با توجه به اینکه آماره F محاسبه شده از $F_{0.05(1,3)} = 10.13$ بزرگتر است، پس رگرسیون معنی‌دار بوده و شیب خط برابر با صفر نیست. نتایج آزمون F را می‌توانید در خروجی برنامه Minitab (جدول ۱۱-۲) نیز مشاهده کنید.

۱۱-۱۱ پیش‌بینی بر اساس مدل رگرسیونی

پس از برآورد مدل مناسب برای ارتباط خطی موجود بین میزان کود و عملکرد، می‌توان از آن در پیش‌بینی استفاده کرد. پیش‌بینی تنها برای مقادیری از متغیر مستقل مجاز است که در دامنه x های بکار رفته در برازش مدل، واقع باشند. زیرا در خارج از این فاصله اعتبار مدل از بین می‌رود و پیش‌بینی بر اساس خط رگرسیونی ممکن است دور از واقعیت باشد. برای مثال، میزان کود بر میزان محصول اثر مستقیم دارد به این معنی که افزایش کود باعث فراوانی محصول می‌شود ولی واضح است که مصرف بیش از حد کود هم کاهش میزان محصول را در پی خواهد داشت. لذا مدل رگرسیونی که بر اساس سطوح مناسب کود ساخته شود پیش‌بینی خوبی برای سطوح بالای کود ارائه نمی‌دهد.

پیش‌بینی می‌تواند به دو صورت باشد. می‌توان به ازای یک مقدار x ، میانگین همه مقادیر y یعنی $E(y)$ را پیش‌بینی نمود یا اینکه به ازای همان مقدار x ، مقدار y را فقط در یک آزمایش منفرد پیش‌بینی کرد. حالت اول در مثال ۱۱-۳ می‌تواند برآورد میانگین درصد جوانه‌زنی در دمای 4°C در چندین تکرار باشد. مثال حالت دوم برآورد درصد جوانه‌زنی در دمای 4°C در یک تکرار است. در هر دو مورد درصد جوانه‌زنی پیش‌بینی شده به صورت

$$\hat{y} = \beta_0 + \beta_1 x = -0.1 + 0.17(4) = 0.58$$

بدست می‌آید. تنها اختلاف بین دو نوع برآورد در دقت آنها می‌باشد. انحراف معیار میانگین جامعه همه مقادیر ممکن y (خطای معیار \hat{y}) به ازای مقدار ثابت $x = x_0$ برابر

$$\sigma_{\hat{y}} = \sigma \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_x}}$$

است. که در آن σ انحراف معیار خطا می‌باشد. این در حالی است که انحراف معیار برآورد \hat{y} (خطای معیار پیش‌بینی) برابر

$$\sigma_{(y-\hat{y})} = \sigma \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_x}}$$

است. فاصله اطمینان $1 - \alpha$ برای میانگین y به ازای $x = x_0$ به صورت زیر به دست می‌آید که در آن $t_{\alpha/2}$ بر اساس $n - 2$ درجه آزادی از جدول t استخراج می‌شود و s نیز برآوردی از σ است که از روی نمونه محاسبه می‌شود. در مثال ۱۱-۳ فاصله اطمینان ۰/۹۵ برای میانگین y ها به ازای $x = 4$ به صورت زیر است:

$$\hat{y} \pm t_{\alpha/2} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_x}} = 2.7 \pm 3.182(0.605) \sqrt{\frac{1}{5} + \frac{(4 - 3)^2}{10}} = 2.7 \pm 1.05$$

که در آن $\alpha/2$ برابر با ۰/۰۲۵ در نظر گرفته شده است. لذا فاصله اطمینان ۰/۹۵ برای میانگین درصد جوانه‌زنی در دمای $4^\circ C$ بصورت ۱/۶۵ تا ۳/۷۵ می‌باشد. این فاصله با استفاده از نمونه موجود به دست آمده است و در صورتی که با استفاده از یک نمونه بزرگتر اطلاعات بیشتری حاصل می‌شد، فاصله اطمینان کوتاه‌تری محاسبه می‌شد.

فاصله پیش‌بینی $1 - \alpha$ برای یک y منفرد به ازای $x = x_0$ به صورت

$$\hat{y} \pm t_{\alpha/2} s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_x}}$$

می‌باشد. در مثال ۱۱-۳ فاصله پیش‌بینی ۰/۹۵ برای یک y منفرد به ازای $x = 4$ به صورت

$$\hat{y} \pm t_{\alpha/2} s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_x}} = 2.7 \pm 3.182(0.605) \sqrt{1 + \frac{1}{5} + \frac{(4 - 3)^2}{10}} = 2.7 \pm 2.19$$

به دست می‌آید. لذا اگر یک بار دیگر دمای $4^\circ C$ بررسی شود، درصد جوانه‌زنی با احتمال ۰/۹۵ بین ۰/۵۱ تا ۴/۸۹ قرار خواهد گرفت. در صورت استفاده از یک نمونه بزرگتر فاصله پیش‌بینی دقیقتری (با طول کوتاه‌تر) به دست خواهد آمد. در جدول ۱۱-۳ فاصله اطمینان برای $E(y)$ و فاصله پیش‌بینی ۰/۹۵ برای y با استفاده از برنامه Minitab برای چند مقدار مختلف x نشان داده شده است.

جدول ۱۱-۳. حل مثال ۱۱-۳ و خروجی مربوطه در برنامه Minitab.

مراحل:					
۱. مشاهدات مربوط به متغیر مستقل (x) و وابسته (y) را در ستون‌های اول و دوم صفحه داده‌ها وارد کنید.					
۲. وارد مسیر Stat > Regression > Regression... شوید:					
۳. در مقابل گزینه Response ستون دوم و در مقابل گزینه Predictors ستون اول (متغیر مستقل) را انتخاب کنید.					
۴. بر روی گزینه Options کلیک کرده و در پنجره مربوطه در مقابل گزینه Predictions intervals for new... ستون مربوط به متغیر مستقل را وارد کنید (البته به جای آن می‌توان یک مقدار جدید x را وارد کرد).					
۴. بر روی گزینه OK کلیک کنید.					

Obs	New Fit	SE Fit	95% CI	95% PI
1	0.600	0.469	(-0.893; 2.093)	(-1.838; 3.038)
2	1.300	0.332	(0.245; 2.355)	(-0.897; 3.497)
3	2.000	0.271	(1.138; 2.862)	(-0.111; 4.111)
4	2.700	0.332	(1.645; 3.755)	(0.503; 4.897)
5	3.400	0.469	(1.907; 4.893)	(0.962; 5.838)

البته باید توجه داشت که نباید از مدل برای پیش‌بینی درصد جوانه‌زنی به ازای مقادیر دمای کمتر از ۱ و بیشتر از ۵ درجه استفاده کرد. زیرا چنین مقادیری خارج از دامنه داده‌هایی هستند که مدل بر اساس آنها برآورد شده است. ممکن است درصد جوانه‌زنی به ازای دمای کمتر از ۱ و بیشتر از ۵ درجه، پراکندگی متفاوتی پیدا کند.

مثال ۱۱-۵. جدول زیر عملکرد گندم را بر حسب مقادیر مختلف کود نشان داده است. به فرض آنکه ۵۵۰ گرم کود مصرف شود، مطلوب است الف) فاصله اطمینان ۹۵٪ میانگین میزان محصول در همه کرت‌ها ب) فاصله اطمینان ۹۵٪ میزان محصول فقط در یک کرت گندم.

x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(y - \bar{y})^2$	$(x - \bar{x})(y - \bar{y})$	\hat{y}	$y - \hat{y}$	$(y - \hat{y})^2$
۱۰۰	۴۰	-۳۰۰	-۲۰	۹۰۰۰۰	۴۰۰	-۶۰۰۰	۴۲/۳	-۲/۳	۵/۲۹
۲۰۰	۵۰	-۲۰۰	-۱۰	۴۰۰۰۰	۱۰۰	-۲۰۰۰	۴۸/۲	۱/۸	۳/۲۴
۳۰۰	۵۰	-۱۰۰	-۱۰	۱۰۰۰۰	۱۰۰	۱۰۰۰	۵۴/۱	-۴/۱	۱۶/۸۱
۴۰۰	۷۰	۰	۱۰	۰	۱۰۰	۰	۶۰/۰	۱۰/۰	۱۰۰/۰
۵۰۰	۶۵	۱۰۰	۵	۱۰۰۰۰	۲۵	۵۰۰	۶۵/۹	-۰/۹	۰/۸۱
۶۰۰	۶۵	۲۰۰	۵	۴۰۰۰۰	۲۵	۱۰۰۰	۷۱/۸	-۶/۸	۳۶/۶۴
۷۰۰	۸۰	۳۰۰	۲۰	۹۰۰۰۰	۴۰۰	۶۰۰۰	۷۷/۷	۲/۳	۵/۲۹

$$\sum x = 2800 \quad \sum y = 420 \quad \sum xy = 184500 \quad \sum x^2 = 1400000 \quad \bar{x} = 400 \quad \bar{y} = 60$$

$$\hat{\beta}_1 = \frac{\sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} = \frac{184500 - \frac{2800 \times 420}{7}}{1400000 - \frac{2800^2}{7}} = \frac{16500}{280000} = 0.059$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 60 - (0.059)400 = 36.43, \quad s^2 = \frac{\sum (y - \hat{y})^2}{n - 2} = 350.8 = 0.96$$

(الف)

$$\hat{y} \pm t_{\alpha/2} s \sqrt{\frac{1}{n} + \frac{(x_i - \bar{x})^2}{SS_x}} = \hat{y} \pm t_{0.025, 25} 0.96 \sqrt{\frac{1}{7} + \frac{(550 - \bar{x})^2}{SS_x}}$$

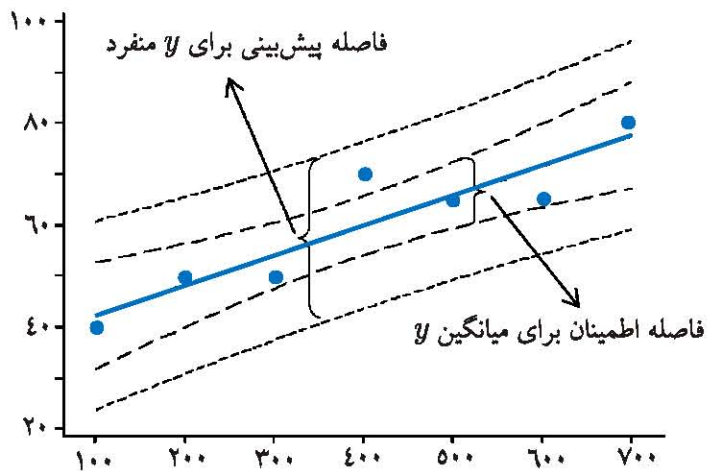
$$38.8 \pm 2.0571(0.96) \sqrt{\frac{1}{7} + \frac{(550 - 400)^2}{280000}} = 38.8 \pm 2.0571(2.816) = 38.8 \pm 7.24$$

(ب)

$$\hat{y} \pm t_{\alpha/2} s \sqrt{1 + \frac{1}{n} + \frac{(x_i - \bar{x})^2}{SS_x}} = \hat{y} \pm t_{0.025, 25} 0.96 \sqrt{1 + \frac{1}{7} + \frac{(550 - \bar{x})^2}{SS_x}}$$

$$38.8 \pm 2.0571(0.96) \sqrt{1 + \frac{1}{7} + \frac{(550 - 400)^2}{280000}} = 38.8 \pm 2.0571(6.09) = 38.8 \pm 16.974$$

طول این فاصله (۱۶/۹۷۴) بیش از دو برابر طول فاصله بدست آمده در قسمت الف (۷/۲۴) است. ملاحظه می شود پیش بینی یک مشاهده منفرد، مشکل تر از پیش بینی یک میانگین است. نمودار ۹-۱۱ فاصله اطمینان را برای میانگین y و فاصله پیش بینی برای y منفرد را نشان می دهد.



نمودار ۹-۱۱. فاصله پیش‌بینی
 ۹۵٪ برای یک y منفرد و فاصله
 اطمینان ۹۵٪ برای میانگین y به
 ازای مقادیر مختلف x در مثال
 ۵-۱۱.

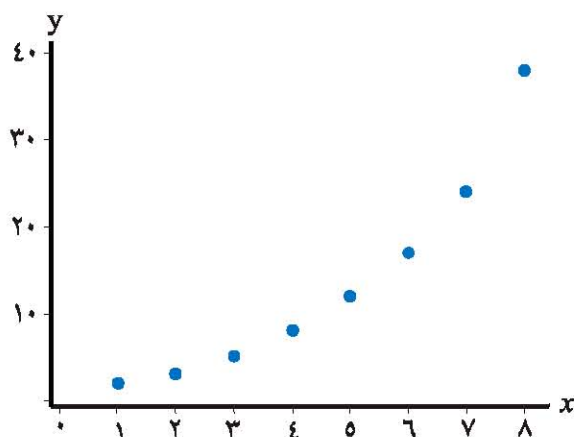
۱۱-۱۲ تبدیل داده‌ها در رگرسیون

گاهی رابطه بین دو متغیر x و y خطی نبوده و نقاط مربوطه در نمودار پراکنش در راستای یک خط مستقیم قرار نمی‌گیرند و لذا مقدار ضریب تعیین (R^2) در رگرسیون خطی کوچک خواهد بود. در چنین مواردی، ممکن است رابطه غیر خطی را بتوان به کمک یک تبدیل به یک رابطه خطی تبدیل کرد. تبدیل داده‌های یک متغیر با توجه به روند تغییرات آن انجام می‌شود. برای مثال اگر مقادیر y در مقایسه به مقادیر x با سرعت بیشتری افزایش یابند، ممکن است x و $\ln y$ یک رابطه خطی با هم داشته باشند یا ممکن است روند افزایش y ها طوری باشد که x و \sqrt{y} یک رابطه خطی ایجاد کند.

مثال ۱۱-۶. فرض کنید میزان فروش یک کالای خاص طی ماه‌های مختلف به صورت زیر بوده است.

x (ماه)	۱	۲	۳	۴	۵	۶	۷	۸
y (میزان فروش)	۲	۳	۵	۸	۱۲	۱۷	۲۴	۳۸

نمودار پراکنش داده‌ها نشان می‌دهد که نقاط در راستای یک خط مستقیم قرار ندارند. لذا رابطه خطی بین دو متغیر وجود ندارد و استفاده از مدل رگرسیون خطی برای بررسی رابطه بین این دو متغیر مناسب نخواهد بود. کم بودن مقدار R^2 نیز گواهی بر این ادعا است (جدول ۱۱-۴).



نمودار ۱۱-۱۰. نمودار پراکنش داده‌های مثال ۱۱-۶.

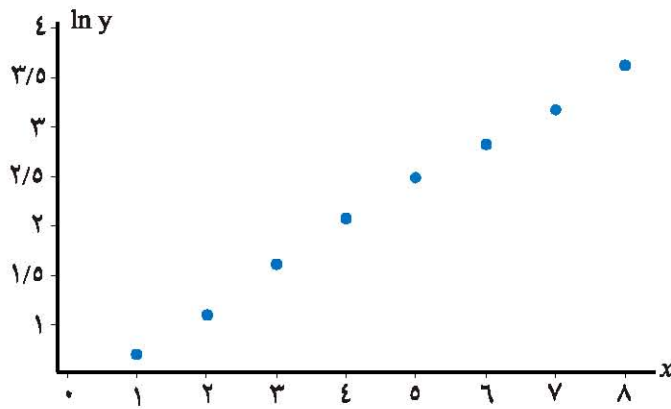
جدول ۱۱-۴. خروجی برنامه Minitab برای مثال ۱۱-۶.

The regression equation is					
$y = - 7.64 + 4.73 x$					
Predictor	Coef	SE Coef	T	P	
Constant	-7.643	3.651	-2.09	0.081	
Y	4.7262	0.7230	6.54	0.001	
S = 4.68555 R-Sq = 87.7% R-Sq(adj) = 85.6%					
Analysis of Variance					
Source	DF	SS	MS	F	P
Regression	1	938.15	938.15	42.73	0.001
Residual Error	6	131.73	21.95		

در اینجا مقادیر y با سرعت بیشتری نسبت به مقادیر x افزایش یافته است. برای خطی کردن چنین رابطه‌ای ممکن است تبدیل y به $\ln y$ بصورت زیر مفید باشد.

x (سال)	۱	۲	۳	۴	۵	۶	۷	۸
$\ln y$	۰/۶۹	۱/۰۲	۱/۶۱	۲/۰۸	۲/۴۸	۲/۸۳	۳/۱۸	۳/۶۴

نمودار پراکنش داده‌ها (نمودار ۱۱-۱۱) نشان می‌دهد که مشاهدات تبدیل یافته در راستای یک خط مستقیم قرار دارند. ضریب تعیین رگرسیون خطی برای این مشاهدات تبدیل شده برابر با ۹۹/۷٪ است (جدول ۱۱-۵)، لذا یک رابطه خطی شدید بین مقادیر x و مقادیر $\ln y$ وجود دارد.



نمودار ۱۱-۱۱. نمودار پراکنش داده‌های مثال ۱۱-۶ پس از تبدیل لگاریتمی.

در این مثال رابطه رگرسیونی غیر خطی با استفاده از تبدیل لگاریتمی یکی از متغیرها به یک رابطه رگرسیونی خطی با ضریب تعیین بالا تبدیل شد. در صورتی که تبدیل لگاریتمی چاره‌ساز نباشد باید تبدیل‌های دیگری را برای داده‌های یکی یا هر دو متغیر آزمود. از تبدیل‌های معمول می‌توان به معکوس کردن داده‌ها ($1/y$ و یا $1/x$)، گرفتن ریشه دوم ($y^{1/2}$ یا $x^{1/2}$) یا ریشه چهارم از داده‌ها ($y^{1/4}$ یا $x^{1/4}$) اشاره کرد. در هر صورت تبدیلی مناسبتر است که به ازای آن، مقدار R^2 رگرسیون خطی برآورد شده به یک نزدیکتر باشد. جدول ۱۱-۶ برخی مدل‌های غیر خطی و تبدیل‌های خطی کننده متناظر آنها را نشان می‌دهد.

جدول ۱۱-۵. خروجی برنامه Minitab برای داده‌های بالا.

The regression equation is: $\ln y = 0.322 + 0.418 x$					
Predictor	Coef	SE Coef	T	P	
Constant	0.32224	0.05045	6.39	0.001	
C1	0.417679	0.009990	41.81	0.000	
S = 0.0647432 R-Sq = 99.7% R-Sq(adj) = 99.6%					
Analysis of Variance					
Source	DF	SS	MS	F	P
Regression	1	7.3272	7.3272	1748.02	0.000
Residual Error	6	0.0252	0.0042		
Total	7	7.3523			

جدول ۱۱-۶. برخی مدل‌های خطی شونده و تبدیل‌های خطی کننده آنها.

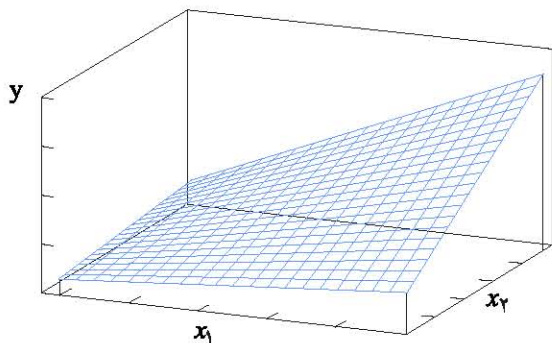
مدل غیر خطی	تبدیل	فرم خطی
$y = \alpha x^\beta$	$y' = \log y, x' = \log x$	$y' = \log \alpha + \beta x'$
$y = \alpha e^{\beta x}$	$y' = \ln y$	$y' = \ln \alpha + \beta x$
$y = \alpha + \beta \log x$	$x' = \log x$	$y = \alpha + \beta x'$
$y = \frac{x}{\alpha x - \beta}$	$y' = \frac{1}{y}, x' = \frac{1}{x}$	$y' = \alpha + \beta x'$

۱۱-۱۳ رگرسیون چندگانه

با توجه به اینکه در عمل غالباً دو یا چند متغیر بر روی متغیر وابسته تأثیر می‌گذارند، بررسی تأثیر همزمان و خطی دو یا چند متغیر مستقل بر متغیر وابسته جالب خواهد بود. برای مثال میزان محصول گندم نه فقط به مقدار کود، بلکه به میزان بارندگی، میزان درجه حرارت، میزان نور و غیره بستگی دارد. در این صورت کمبود نور یا حرارت، عملکرد را محدود می‌کند حتی اگر کود کافی نیز در دسترس گیاه قرار گرفته باشد. اگر متغیر وابسته بصورت خطی تحت تأثیر k متغیر مستقل باشد، معادله رگرسیون خطی چند متغیره را می‌توان بصورت

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon$$

نوشت که در آن پارامترهای مدل و ε جمله خطا را نشان می‌دهد. این معادله تعمیمی از حالت خطی ساده است. معادله $E(y) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$ که معادله میانگین پاسخ یا رویه پاسخ نامیده می‌شود در حالت کلی نشان‌دهنده یک رویه در فضای $k + 1$ است. در حالت خاص که $k = 2$ باشد، معادله $E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$ نشان‌دهنده صفحه‌ای در فضای سه بعدی است (نمودار ۱۱-۱۲).



نمودار ۱۱-۱۲. صفحه میانگین پاسخ به

ازای $k = 2$ یعنی

$$.E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

برآورد مدل رگرسیونی در حالت چند متغیره همانند حالت تک متغیره با استفاده از روش کمترین توان‌های دوم و با مینیمم کردن مجموع توان‌های دوم خطا صورت می‌پذیرد. در حالت دو متغیره می‌توان نوشت

$$SS_E = \sum (y_i - \hat{y})^2 = \sum [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2})]^2$$

اگر مشتق‌های جزئی نسبت به $\hat{\beta}_0$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ محاسبه و برابر صفر قرار داده شوند، برآوردگرهای ضرایب رگرسیونی از حل همزمان معادلات

$$n\hat{\beta}_0 + \hat{\beta}_1 \sum x_1 + \hat{\beta}_2 \sum x_2 = \sum y$$

$$\hat{\beta}_0 \sum x_1 + \hat{\beta}_1 \sum x_1^2 + \hat{\beta}_2 \sum x_1 x_2 = \sum x_1 y$$

$$\hat{\beta}_0 \sum x_2 + \hat{\beta}_1 \sum x_1 x_2 + \hat{\beta}_2 \sum x_2^2 = \sum x_2 y$$

بدست می‌آیند. در واقع روش محاسبه ضرایب مشابه روش رگرسیون ساده است با این تفاوت که در اینجا محاسبات پیچیده‌تر خواهند بود، به ویژه اگر بیش از دو متغیر در مدل وجود داشته باشد. به همین دلیل در اینجا به بررسی رگرسیون چندگانه با استفاده از برنامه آماری Minitab پرداخته می‌شود.

مثال ۱۱-۷. میزان عملکرد ۱۵ کرت از یک نوع گیاه بر اساس میزان کود (x_1) و تعداد دفعات آبیاری (x_2) در جدول زیر آورده شده است. معادله رگرسیون را برآورد کنید به طوری که بتوان با داشتن میزان کود و تعداد دفعات آبیاری، عملکرد را پیش‌بینی کرد.

میزان عملکرد (y)	کود (x_1)	دفعات آبیاری (x_2)
۲۸۴۱	۵۰	۱۱
۱۸۷۶	۲۱	۸
۲۹۳۴	۳۸	۱۰
۱۵۵۲	۱۸	۱۰
۳۰۶۵	۴۳	۱۲
۳۶۷۰	۶۵	۱۲
۲۰۰۵	۵۰	۵
۳۲۱۵	۴۸	۸
۱۹۳۰	۱۷	۸
۲۰۱۰	۷۰	۶
۳۱۱۱	۲۰	۹
۲۸۸۲	۲۹	۹
۱۶۸۳	۱۵	۵
۱۸۱۷	۱۴	۷
۴۰۶۶	۶۰	۱۳

معادله رگرسیون برآزش شده همانگونه که در خروجی Minitab نشان داده شده است $\hat{y} = 277 + 15/3x_1 + 195x_2$ می‌باشد. در این مثال $\hat{\beta}_1$ برابر با $15/3$ و $\hat{\beta}_2$ برابر با 195 به دست آمده است. این بدان معنی است که با ثابت بودن تعداد دفعات آبیاری، هر واحد کود عملکرد را به اندازه $15/3$ واحد و با ثابت بودن میزان کود، هر بار آبیاری عملکرد را به اندازه 195 واحد افزایش می‌دهد. همچنین $\hat{\beta}_0$ برابر با 277 به دست آمده است.

در قسمت بعدی خطای معیار ($SE\ Coef$) برای هر کدام از این ضرایب محاسبه شده است و به کمک آنها می‌توان فاصله اطمینان $1 - \alpha$ را برای هر کدام از ضرایب به دست آورد. برای مثال فاصله اطمینان 95% برای β_1 بر اساس رابطه $\hat{\beta}_1 \pm t_{\alpha/2} S_e$ برابر با $2/179(6/97) \pm 15/3$ است. به عبارت دیگر با احتمال 95% مقدار β_1 در فاصله $0/11$ تا $30/48$ قرار دارد. با توجه به اینکه این فاصله صفر را در بر نمی‌گیرد، لذا β_1 در سطح 5% اختلاف معنی‌داری با صفر دارد.

جدول ۷-۱۱. حل مثال ۷-۱۱ و خروجی مربوطه در برنامه Minitab.

مراحل:

۱. داده‌های متغیر x_1 ، x_2 و y را در ستون‌های جداگانه وارد کنید.

Stat > Regression > Regression...		۲. وارد مسیر مقابل شوید:			
در مقابل گزینه Response ستون y و در مقابل گزینه Predictors متغیرهای مستقل را انتخاب کنید.		۳.			
		۴. بر روی OK کلیک کنید.			
The regression equation is Y = 277 + 15.3 X1 + 195 X2					
Predictor	Coef	SE Coef	T	P	
Constant	276.8	485.5	0.57	0.579	
X1	15.262	6.970	2.19	0.049	
X2	195.40	54.07	3.61	0.004	
S = 481.650 R-Sq = 67.7% R-Sq(adj) = 62.3%					
Analysis of Variance					
Source	DF	SS	MS	F	P
Regression	2	5834582	2917291	12.58	0.001
Residual Error	12	2783846	231987		
Total	14	8618428			

در قسمت بعدی خروجی، انحراف معیار خطا ($S = 481.65$) و همچنین ضریب تعیین ($R^2 = 67.7\%$) نشان داده شده است. انحراف معیار خطا میزان پراکندگی داده‌ها را حول رگرسیون نشان می‌دهد و در حالت کلی به صورت

$$S = \sqrt{\frac{SS_E}{df_E}} = \sqrt{\frac{\sum (y_i - \hat{y})^2}{n - k - 1}}$$

محاسبه می‌شود که در آن، n برابر با تعداد مشاهدات و k تعداد متغیرهای مستقل است. درجه آزادی در این مثال ۱۲ می‌باشد. ضریب تعیین (R^2) نیز که از تقسیم SS رگرسیون بر SS کل ($SS_y = SS_{Total}$) به دست می‌آید بیانگر سهمی از تغییرات متغیر وابسته است که توسط رگرسیون توجیه می‌شود. مقدار R^2 در اینجا برابر با 67.7% است. با توجه به اینکه R^2 خیلی به ۱ نزدیک نیست پیش‌بینی عملکرد بر اساس تعداد آبیاری و میزان کود با استفاده از این مدل خیلی دقیق نخواهد بود. برنامه‌های آماری همچنین ضریب تعیین تصحیح شده یا R^2_{adj} را نیز به دست می‌دهند که از رابطه

$$R^2_{adj} = 1 - \frac{SS_E / (n - k - 1)}{SS_y / (n - 1)}$$

بدست می‌آید. بعضی از محققین استفاده از آماره R_{adj}^y را در مورد رگرسیون‌های چندگانه ترجیح می‌دهند، زیرا ضریب تعیین معمولی (R^y) همواره با اضافه شدن یک متغیر جدید به مدل رگرسیون افزایش می‌یابد یا حداقل کاهش نمی‌یابد. ضریب تعیین تصحیح شده با استفاده از میانگین توان‌های دوم به جای مجموع توان‌های دوم این مشکل را حل می‌کند. اگر تفاوت R^y با R_{adj}^y زیاد باشد، متغیر یا متغیرهایی که مشارکت معنی‌داری در برازش ندارند در مدل رگرسیونی گنجانده شده‌اند. از این خاصیت R_{adj}^y ممکن است حتی برای گزینش متغیرهای مستقل مفید در مدل‌سازی استفاده شود.

تحلیل واریانس در خروجی Minitab (جدول ۱۱-۷) نیز نشان می‌دهد که تغییرات کل به دو قسمت مجموع توان‌های دوم خطا و مجموع توان‌های دوم رگرسیون تفکیک شده است.

۱۱-۱۴ رگرسیون چندجمله‌ای

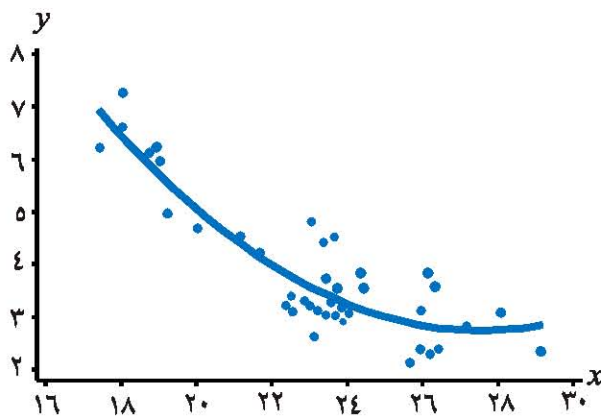
گاهی در رگرسیون ساده مشاهده می‌شود که نمودار پراکنش داده‌ها منحنی است و نمی‌توان تبدیل مناسبی برای خطی کردن رابطه به کار برد. در چنین مواردی متغیر وابسته y تابعی از درجات بالاتر x نیز می‌باشد. به اینگونه مدل‌ها رگرسیون چندجمله‌ای گفته می‌شود. برای مثال مدل چندجمله‌ای $y = \beta_0 + \beta_1x + \beta_2x^2$ رگرسیون درجه دو نامیده می‌شود و برای حالاتی مناسب است که با وجود یک متغیر مستقل، نمودار پراکنش داده‌ها ابتدا یک کاهش و سپس افزایش تدریجی یا ابتدا یک افزایش سپس کاهش تدریجی نشان می‌دهند. به مدل چندجمله‌ای $y = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3$ نیز رگرسیون درجه سه گفته می‌شود و برای مواردی مناسب است که با وجود یک متغیر مستقل، نمودار پراکنش داده‌ها دو نقطه ماکسیمم و یک نقطه مینیمم یا دو نقطه مینیمم و یک نقطه ماکسیمم را نشان دهد.

تحلیل رگرسیون چندگانه همانند رگرسیون چندجمله‌ای است با این تفاوت که در اینجا یک متغیر مستقل وجود دارد. اگر داده‌های متغیر مستقل در ستون جدیدی به توان دو برسند و اثر x^2 و بر متغیر وابسته در قالب رگرسیون چندگانه بررسی شود، مدل رگرسیون دوجمله‌ای ایجاد شده است.

مثال ۱۱-۸ در تحقیقی به منظور بررسی اثر میانگین حداقل دمای روزانه در یک ماه اول بعد از کاشت بر روی عملکرد برنج بر حسب تن در هکتار، داده‌های زیر به دست آمد.

عملکرد (y)	میانگین دما (x)	عملکرد (y)	میانگین دما (x)
۲/۳	۲۹/۲	۳	۲۳/۴
۳/۱	۲۸/۱	۴/۴	۲۳/۴
۲/۸	۲۷/۲	۳/۱	۲۳/۲
۲/۴	۲۶/۴	۲/۶	۲۳/۱
۳/۶	۲۶/۳	۴/۸	۲۳/۱
۲/۳	۲۶/۲	۳/۲	۲۳/۰
۳/۸	۲۶/۲	۳/۳	۲۲/۹
۳/۱	۲۶/۰	۳/۱	۲۲/۵
۲/۴	۲۵/۹	۳/۴	۲۲/۵
۲/۱	۲۵/۷	۳/۲	۲۲/۴
۳/۵	۲۴/۵	۴/۲	۲۱/۷
۳/۸	۲۴/۴	۴/۵	۲۱/۲
۳/۱	۲۴/۰	۴/۷	۲۰/۰
۲/۹	۲۳/۹	۵/۰	۱۹/۲
۳/۲	۲۳/۹	۶/۲	۱۹/۰
۳/۰	۲۳/۷	۶/۰	۱۹/۰
۴/۵	۲۳/۷	۶/۱	۱۸/۸
۳/۵	۲۳/۷	۷/۳	۱۸/۰
۳/۳	۲۳/۶	۶/۶	۱۸/۰
۳/۷	۲۳/۵	۶/۲	۱۷/۴

نمودار پراکنش نشان می‌دهد که داده‌ها در راستای یک منحنی طوری قرار گرفته‌اند که ممکن است رگرسیون درجه دو برازش خوبی بر نقاط داشته باشد.



نمودار ۱۱-۱۳. نمودار پراکنش داده‌های مربوط به عملکرد برنج در مقابل میانگین حداقل دمای روزانه و منحنی برازش داده شده بر آنها.

جدول ۸-۱۱ برآزش منحنی رگرسیون درجه دو به داده‌های مثال ۸-۱۱

مراحل:					
۱. مشاهدات متغیر y و x را بترتیب در ستون‌های اول و دوم صفحه داده‌ها وارد کنید.					
۲. مربع داده‌های ستون مربوط به متغیر x را تحت عنوان x_2 در ستون سوم وارد کنید. برای این منظور می‌توان از Calculator... موجود در منوی Calc استفاده کرد.					
۳. وارد مسیر مقابل شوید: Stat > Regression > Regression...					
۴. در مقابل گزینه Response ستون y و در مقابل گزینه Predictors متغیرهای ستون‌های x و x_2 را انتخاب کنید.					
۵. بر روی OK کلیک کنید.					
The regression equation is $y = 33.7 - 2.24 x + 0.0408 x^2$					
Predictor	Coef	SE Coef	T	P	
Constant	33.695	4.876	6.91	0.000	
x	-2.2469	0.4299	-5.23	0.000	
x2	0.040764	0.009391	4.34	0.000	
S = 0.587621 R-Sq = 80.8% R-Sq(adj) = 79.8%					
Analysis of Variance					
Source	DF	SS	MS	F	P
Regression	2	53.932	26.966	78.09	0.000
Residual Error	37	12.776	0.345		
Total	39	66.708			

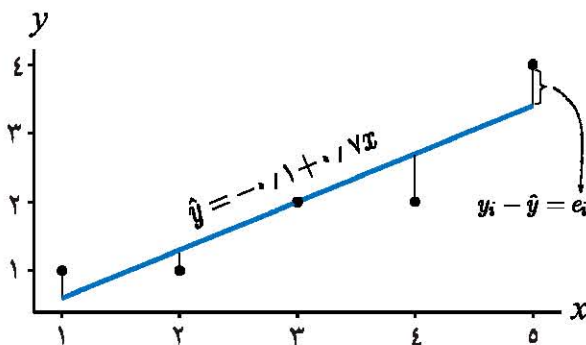
همانگونه که در خروجی Minitab نشان داده شده است معادله منحنی رگرسیون برآزش شده بصورت $\hat{y} = 33.7 - 2.24x + 0.04x^2$ می‌باشد. اگر به جای رگرسیون درجه دو، رگرسیون خطی ساده بر داده‌های متغیر وابسته و مستقل برآزش داده می‌شد، ضریب تعیین برابر با $0.71/1$ به دست می‌آمد، که خیلی کمتر از مقدار فعلی $(0.80/1)$ است.

۱۱-۱۵ بازبینی مدل آماری

به طور کلی در تحلیل رگرسیونی، مدلی بر اساس کمترین مجموع توان‌های دوم خطا برآزش داده شده و بعضاً فرض‌های مختلفی مورد آزمون قرار می‌گیرد. غالباً لازم است که مدل رگرسیونی برآورد شده به طور دقیقتری بازبینی شود. این بازبینی با اهداف مختلفی انجام

می‌شود. مثلاً اینکه آیا می‌توان مدل مناسبتری که دارای خطای کمتر و دقت بیشتری باشد برآورد نمود؟ یکی از معمولترین کارها برای کشف عیوب مدل و ارائه راه حلی مناسب برای اصلاح آن، مشاهده نمودار پراکنش مانده‌ها می‌باشد. مانده‌ها در واقع انحرافات نقاط از رگرسیون برازش داده شده هستند و نمی‌توانند پس از برازش مدل بر داده‌ها تبیین شوند. برای مثال اگر عملکرد (متغیر وابسته) در مقادیر کودی ۱، ۲، ۳، ۴ و ۵ به ترتیب ۱، ۲، ۲، ۴ و ۴ باشد، مانده‌ها یا خطاها را می‌توان پس از برازش خط رگرسیون به صورت زیر به دست آورد.

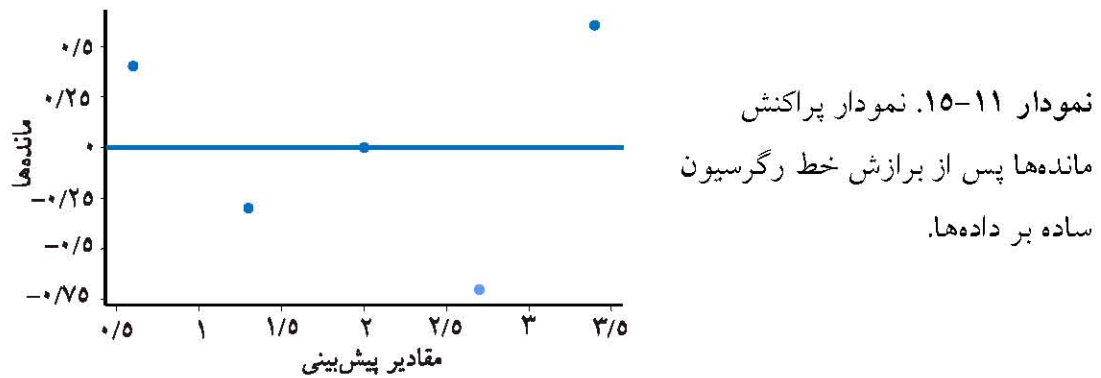
$$e_i = y_i - \hat{y}$$



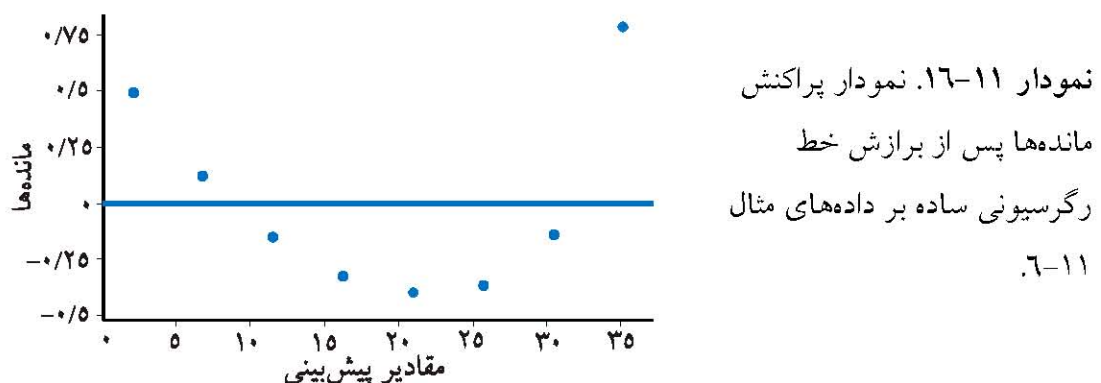
نمودار ۱۱-۱۴. مانده‌ها (e_i ها) پس از برازش خط رگرسیون بر مشاهدات.

مهمترین فرض در هر مدلی این است که مانده‌ها مستقل از هم و دارای توزیع نرمال باشند. سایر فرض‌های رگرسیون در بند ۱۱-۶ آورده شده‌اند. در واقع همه این فرض‌ها دلالت بر این دارند که مانده‌ها نباید از هیچ الگوی مشخصی پیروی کرده و باید به صورت تصادفی مشاهده شوند. مشاهده هر گونه الگوی منظمی در مورد مانده‌ها برقراری فرض‌های رگرسیون را زیر سؤال می‌برد. روش معمول در بررسی مانده‌ها مشاهده نمودار پراکنش مانده‌ها در مقابل مقادیر پیش‌بینی شده می‌باشد.

در برنامه Minitab پس از رفتن به مسیر **Stat > Regression > Regression...** برای تجزیه رگرسیون، با کلیک روی گزینه **Graphs...** و انتخاب **Residuals versus fits** می‌توان نمودار پراکنش مانده‌ها را در مقابل مقادیر برآورد شده مشاهده نمود. نمودار پراکنش مانده‌ها برای درصد جوانه‌زنی به ازای دماهای مختلف پس از برازاندن خط رگرسیون ساده، در نمودار ۱۱-۱۵ نشان داده شده و هیچگونه روندی در آن مشاهده نمی‌شود.



حال اگر به نمودار پراکنش مانده‌های مربوط به مثال ۱۱-۶ پس از برازاندن خط رگرسیون ساده نگاه کنید می‌بینید که به جای اینکه مانده‌ها حول محور \hat{y} به صورت تصادفی پراکنده شده باشند، ابتدا یک کاهش سپس افزایش تدریجی دارند. در چنین مواردی مدل خطی برازش داده‌شده مناسب نبوده و باید حالت‌های دیگری از جمله مدل‌های غیر خطی بررسی شوند. همچنین می‌توان تبدیل مناسبی را در مورد داده‌های یکی یا هر دو متغیر اعمال کرده و سپس رگرسیون خطی ساده را در مورد داده‌های تبدیل شده به کار برد. به یاد داریم که این مشکل در مثال مذکور از طریق تبدیل لگاریتمی متغیر وابسته حل شد.



۱۱-۱. داده‌های زیر را در نظر بگیرید. الف) داده‌ها را بر روی محور مختصات نشان دهید. ب) مقادیر r و r^2 را محاسبه کنید. ج) آیا شواهد کافی مبنی بر وجود رابطه خطی بین x و y وجود دارد؟ در سطح معنی‌داری $\alpha = 0.01$ بررسی کنید.

x	۳	۵	۶	۴	۳	۷	۶	۵	۴	۷
y	۴	۳	۲	۱	۲	۳	۳	۵	۴	۲

۱۱-۲. در تحقیقی به منظور بررسی وجود ارتباط بین قطر درختان در فاصله ۱/۵ متری از زمین و ارتفاع درختان، ۱۲ درخت به طور تصادفی انتخاب و داده‌های زیر به دست آمد.

قطر درخت (cm)	۹/۸	۷/۲	۲۴/۰	۷/۸	۲/۵	۴۰/۲	۱۵/۵	۶۴/۴	۳۱/۵	۱۰/۸	۳/۰	۸/۳
ارتفاع درخت (m)	۱۰/۹	۹/۸	۲۰/۵	۱۱/۰	۵/۵	۲۰/۴	۱۷/۵	۲۶/۸	۲۵/۶	۱۲/۳	۷/۹	۱۳/۶

الف) مقدار r^2 را محاسبه و تفسیر کنید. ب) معادله خط رگرسیونی ارتفاع بر روی قطر درختان را با روش حداقل مجموع توان‌های دوم به دست آورده و پارامترهای برآورد شده را تفسیر نمایید. ج) مقادیر برازش داده شده و خطاها را برای هر یک از نقاط به دست آورده و آنها را بر روی محور مختصات نشان دهید. آیا نقاط در راستای یک خط راست قرار می‌گیرند؟ د) واریانس کل، خطا و رگرسیون را به دست آورده معنی‌دار بودن ارتباط خطی را از طریق آزمون F بررسی کنید. ه) انحراف معیار $\hat{\beta}_0$ و $\hat{\beta}_1$ را به دست آورده و با استفاده از آنها فاصله اطمینان ۹۵٪ را برای عرض از مبدأ و شیب خط جامعه پیدا کرده و نتایج را توضیح دهید. و) قطر یک درخت جدید اندازه‌گیری و مقدار آن ۱۰ سانتی‌متر بوده است. مقدار ارتفاع آن را پیش‌بینی کرده و فاصله پیش‌بینی ۹۵٪ را برای این ارتفاع به دست آورید. ز) فرض‌های رگرسیون را با استفاده از خط رگرسیون و نقاط مشاهده شده بر روی محور مختصات مورد بحث قرار دهید.

۱۱-۳. مقدار pH در ۱۷ نوع خاک با درصد رس (x_1)، درصد مواد آلی (x_2) و درصد کربنات کلسیم (x_3) متفاوت به صورت زیر بوده است. معادله رگرسیون را با استفاده از یک برنامه آماری برآورد کنید به طوری که بتوان با داشتن مقدار رس، مقدار مواد آلی و مقدار کربنات کلسیم، pH خاک را پیش‌بینی کرد. موارد مشاهده شده در خروجی برنامه را تفسیر نمایید.

درصد کربنات کلسیم (x_2)	درصد مواد آلی (x_3)	درصد رس (x_1)	pH خاک (y)
۶/۱	۴/۳	۵۱/۱	۷/۱
۰/۰	۲/۶	۲۲/۰	۵/۴
۲/۰	۳/۰	۱۷/۰	۷/۰
۰/۰	۳/۰	۱۶/۸	۶/۱
۰/۰	۴/۰	۵/۵	۳/۷
۰/۱	۳/۳	۲۱/۲	۷/۰
۱۶/۸	۳/۷	۱۴/۱	۷/۴
۱۷/۳	۰/۷	۱۶/۶	۷/۴
۱۵/۶	۳/۷	۳۵/۹	۷/۳
۱۱/۹	۳/۳	۲۹/۹	۷/۵
۲/۸	۳/۱	۲/۴	۷/۴
۶/۲	۲/۸	۱/۶	۷/۴
۰/۳	۱/۸	۱۷/۰	۷/۳
۹/۱	۲/۳	۳۲/۶	۷/۳
۰/۰	۴/۰	۱۰/۵	۴/۰
۲۶/۰	۵/۱	۳۳/۰	۷/۱
۰/۰	۱/۹	۲۶/۰	۵/۶