

رگرسیون خطی ساده: فرض کنیم x و y دو متغیر کمی باشند و $(x_1, y_1), (x_2, y_2), \dots$

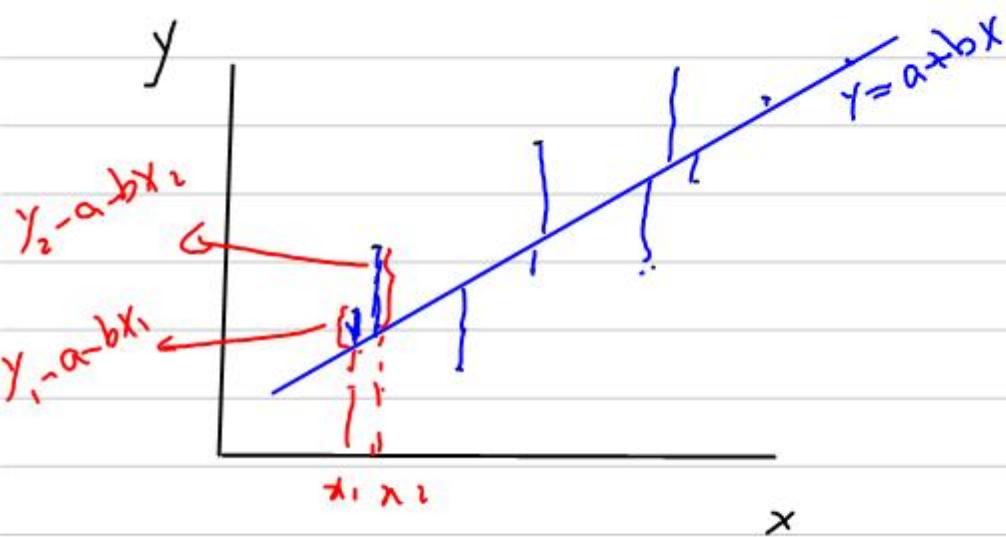
(x_n, y_n) ... یک نمونه تصادفی از یک توزیع متغیر (x, y) باشد

اگر فرض کنیم بین x و y یک رابطه مستقیم وجود دارد، رابطه بین x و y را به صورت $y = a + bx$ می‌نویسند

معادله مستقیم $y = a + bx$ به دست آورد این معادله را خط رگرسیونی می‌گویند

لازم است x متغیر مستقل و y متغیر وابسته باشد و x را متغیر مستقل می‌گویند

این روش را "روش کمترین مربعات" می‌گویند.



روش کمترین مربعات خطی

هدف از این روش آن است که مجموع مربعات خطای را به حداقل برسانیم

$$Q = \sum_{i=1}^n (y_i - a - bx_i)^2 = \min$$

یعنی با تغییر a و b رابطه را به بهترین حالت می‌رسانیم

پس باید

$$\frac{\partial Q}{\partial a} = 0, \quad \frac{\partial Q}{\partial b} = 0$$

$$\frac{\partial Q}{\partial a} = \sum -2(y_i - a - bx_i) = 0 \quad \left\{ \begin{array}{l} na + b \sum x_i = \sum y_i \\ a \sum x_i + b \sum x_i^2 = \sum x_i y_i \end{array} \right.$$

$$\frac{\partial Q}{\partial b} = \sum -2x_i(y_i - a - bx_i) = 0$$

پس از حل دستگاه معادلات

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

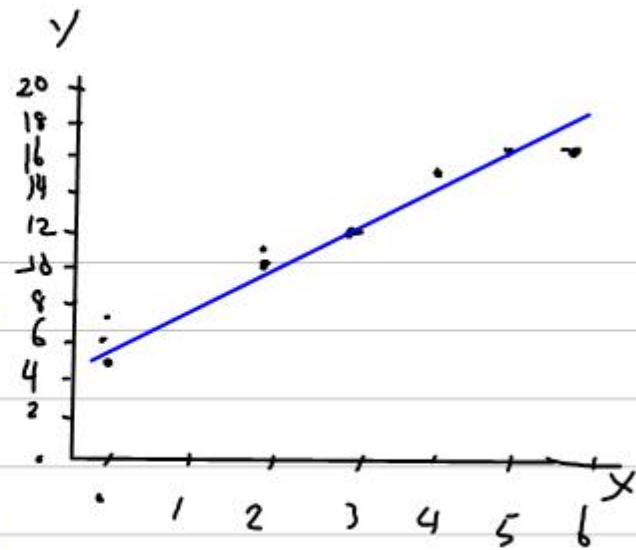
$$a = \bar{y} - b \bar{x}$$

مثال: تعداد ساعت مطالعه و نمره در میان 10 دانشجو به شرح زیر بوده است

| | | | | | | | | | | |
|----------|---|----|----|---|----|----|----|----|----|---|
| x مطالعه | 0 | 2 | 3 | 0 | 2 | 4 | 5 | 6 | 2 | 0 |
| y نمره | 5 | 10 | 12 | 6 | 11 | 15 | 16 | 16 | 10 | 7 |

- الف- ابتدا آرایش داده شده را مرتب کنیم
- ب- معادله خط را بر حسب نمره و ساعت مطالعه رابطه می‌نویسیم
- ج- نمره را می‌توانیم در 2.5 ساعت مطالعه برده‌ایم را تقریباً بدانیم

| | | | | | | | | | | |
|-----------|----|-----|-----|----|-----|-----|-----|-----|-----|----|
| عدد | 0 | 2 | 3 | 0 | 2 | 4 | 5 | 6 | 2 | 0 |
| of x | 5 | 10 | 12 | 6 | 11 | 15 | 16 | 16 | 10 | 7 |
| x_i^2 | 0 | 4 | 9 | 0 | 4 | 16 | 25 | 36 | 4 | 0 |
| $x_i y_i$ | 0 | 20 | 36 | 0 | 22 | 60 | 80 | 96 | 20 | 0 |
| y_i^2 | 25 | 100 | 144 | 36 | 121 | 225 | 256 | 256 | 100 | 49 |



$$\sum x_i = 24 \quad \sum y_i = 108 \quad \sum x_i^2 = 98 \quad \sum x_i y_i = 334 \quad \sum y_i^2 = 1312$$

$$b = \frac{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} = \frac{334 - \frac{24(108)}{10}}{98 - \frac{(24)^2}{10}} = \frac{74.8}{40.4} = 1.85$$

$$a = \bar{y} - b \bar{x} = \frac{108}{10} - 1.85 \left(\frac{24}{10} \right) = 6.36$$

$$y = 6.36 + 1.85x$$

عمرت از جبر

بسیار

ص، احوال، ابر، عبارت از

6.36: میانگین نمره ای که در مطالعه می نوز

1.85: برای هر یک واحد مطالعه نمره این نمره

$$x = 2.5 \Rightarrow y = 6.36 + 1.85(2.5) = 10.985$$

$$r(x, y) = \frac{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}{\sqrt{\sum x_i^2 - \frac{(\sum x_i)^2}{n}}} \sqrt{\sum y_i^2 - \frac{(\sum y_i)^2}{n}} = \frac{334 - \frac{24(108)}{10}}{\sqrt{98 - \frac{(24)^2}{10}}} \sqrt{1312 - \frac{(108)^2}{10}} = 0.975$$

شاخصه‌ها بر اثر مدل

① ضریب تعیین

$$R^2 = r^2$$

ضریب تعیین یعنی چقدر در تغییرات متغیر وابسته توسط متغیر مستقل تبیین می‌شود.

در مثال فوق

$$R^2 = (0.975)^2 = 0.95$$

یعنی 95٪ تغییرات متغیر وابسته توسط متغیر مستقل تبیین می‌شود و فقط 5٪ آن توسط عوامل دیگر تبیین می‌شود.

② ضریب تعیین مدل

| | | |
|---|--------|--------------|
| { | $H_0:$ | مدل خطی نیست |
| | $H_1:$ | مدل خطی هست |

برای انجام آزمون F، بر اساس این آزمون: $\text{تغییرات کل} = \text{تغییرات تبیین شده} + \text{تغییرات تبیین نشده}$

$$SSTO = SSR + SSE$$

$$R^2 = \frac{SSR}{SSTO} = r^2$$

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

$$SSR = a \sum y_i + b \sum x_i y_i - \frac{(\sum y_i)^2}{n}$$

$$SSE = SSTO - SSR$$

پس جدول آزمون Anova به صورت زیر تکمیل می شود

| منبع تغییرات | SS | df | MS | F* |
|--------------|------|-----|-------------------------|-------------------|
| انگیزه | SSR | 1 | $MSR = \frac{SSR}{1}$ | $\frac{MSR}{MSE}$ |
| مانده | SSE | n-2 | $MSE = \frac{SSE}{n-2}$ | |
| کل | SSTO | n-1 | | |

H_0 : ادعای خطای α رد نمی شود
 $F^* > F_{\alpha}(1, n-2)$

در مثال فوق، خطای اولی را در سطح خطای 5٪ در نظر می گیریم

- H_0 : مدل خطای نیست
- H_1 : مدل خطای است

$$SSTO = \sum y_i^2 - \frac{(\sum y_i)^2}{n} = 1312 - \frac{(108)^2}{10} = 145.6$$

$$SSR = a \sum y_i + b \sum x_i y_i - \frac{(\sum y_i)^2}{n} = 6.36(108) + 1.85(334) - \frac{(108)^2}{10} = 138.38$$

$$R^2 = \frac{138.38}{145.6} = 0.95$$

$$SSE = 145.6 - 138.38 = 7.22$$

| منبع تغییرات | SS | df | MS | F* |
|--------------|--------|----|--------|--------|
| انگیزه | 138.38 | 1 | 138.38 | 153.33 |
| مانده | 7.22 | 8 | 0.9025 | |
| کل | 145.6 | 9 | | |

$F^* = 153.33 > F_{0.05}(1, 8) = 5.32$ بنابراین
 H_0 رد می شود پس مدل خطای است

مثال: محقق قصد دارد عملکرد آرزوین صنعت را با استفاده از نمره‌های حقوقی در

در دوره‌ها آرزوین ضمن خدمت و پیش‌بینی نماید. بر این اساس، عددهای زیر از میان آن‌ها به ترتیب در دوره‌ها ۱ تا ۷ در نظر گرفته شده است. آن‌ها را در جدول زیر ثبت کرده است.

| | | | | | | | |
|------------|----|----|----|----|----|----|----|
| نمره حقوقی | 3 | 5 | 2 | 0 | 7 | 1 | 6 |
| آرزوین | 15 | 17 | 15 | 10 | 28 | 10 | 27 |

- الف) معادله خط رگرسیونی برای عملکرد بر حسب تعداد دوره‌ها را تعیین کنید.
 ب) عملکرد معکوس را ۴ دوره‌ها را اندازه‌گیری کنید.
 ج) ضرایب مدل را تفسیر کنید.
 د) ضریب تعیین را محاسبه و تفسیر کنید.
 ه) ضریب همبستگی مدل را در سطح اطمینان ۵٪ بیان کنید.

$a = 8.41$
 $b = 2.63$

| x_i | y_i | $x_i y_i$ | x_i^2 | y_i^2 |
|-------|-------|-----------|---------|---------|
| 3 | 15 | 45 | 9 | 225 |
| 5 | 17 | 85 | 25 | 289 |
| 2 | 15 | 30 | 4 | 225 |
| 0 | 10 | 0 | 0 | 100 |
| 7 | 28 | 196 | 49 | 784 |
| 1 | 10 | 10 | 1 | 100 |
| 6 | 27 | 162 | 36 | 729 |
| 24 | 122 | 528 | 124 | 2452 |

$$b = \frac{528 - \frac{24(122)}{7}}{124 - \frac{(24)^2}{7}} = \frac{109.7}{41.7} = 2.63$$

$$a = \frac{122}{7} - 2.63 \left(\frac{24}{7} \right) = 8.41$$

$$Y = 8.41 + 2.63 X$$

الف مدل =

$$X=4 \Rightarrow Y = 8.41 + 2.63(4) = 18.93$$

ب -

ج - 8.41 : متوسط نمره عمریات کار جمع (دوره اول) نمره نمره

2.63 : به ازای هر دوره، نمره عمریات افزایش می‌کند

->

$$\begin{cases} H_0: & \text{مدل خطی نیست} \\ H_1: & \text{مدل خطی است} \end{cases}$$

$$SSTO = \sum y_i^2 - \frac{(\sum y_i)^2}{n} = 2452 - \frac{(122)^2}{7} = 325.71$$

$$SSR = a \sum y_i + b \sum x_i y_i - \frac{(\sum y_i)^2}{n} = 8.41(122) + 2.63(528) - \frac{(122)^2}{7} = 288.37$$

$$SSE = 325.71 - 288.37 = 37.34$$

$$R^2 = \frac{SSR}{SSTO} = \frac{288.37}{325.71} = 0.89$$

یعنی 89٪ تغییر نمره عمریات توسط مدل دوره‌ها تفسیر می‌شود. / آن را می‌توان در این زمینه در نظر گرفت.

| منبع تغییرات | SS | df | MS | F* |
|--------------|--------|----|--------|------|
| رگرسیون | 288.37 | 1 | 288.37 | 38.6 |
| مانده | 37.34 | 5 | 7.47 | |
| کل | 325.71 | 6 | | |

چون $F^* = 38.6 > F_{0.05} [1, 5] = 6.61$ پس H_0 رد می‌شود یعنی مدل خطی است.

رگرسیون خطی چندگانه

اگر یک متغیر وابسته Y را تعداد K متغیر مستقل X_1, X_2, \dots, X_K در پیش بینی کنیم، خواهیم داشت:

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_K X_K$$

یک رابطه خطی در صورت

به دست آوریم، گوئیم رگرسیون خطی چندگانه داریم.

عرض کنیم $(X_{11}, X_{21}, Y_1), (X_{12}, X_{22}, Y_2), \dots, (X_{1n}, X_{2n}, Y_n)$

یک نمونه نهای از این متغیر X_1, X_2, Y یا نمونه برای n مشاهده صد می باشد.

$$Y = b_0 + b_1 X_1 + b_2 X_2$$

به صورت زیر عمل می‌کنیم -

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$$

نویس Y

ماتریس X با n سطر و 3 ستون به صورت زیر درج می‌شود:

$$X = \begin{bmatrix} 1 & X_{11} & X_{21} \\ 1 & X_{12} & X_{22} \\ \vdots & \vdots & \vdots \\ 1 & X_{1n} & X_{2n} \end{bmatrix}$$

$$Y = XB$$

$$B = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix}$$

ماتریس B

برای حل بردش تیرین مجموع برش خط تیریس B از زیر لایه برش می آید

$$B = (X'X)^{-1} X'Y$$

$$X = \begin{bmatrix} 1 & x_{11} & x_{21} \\ 1 & x_{12} & x_{22} \\ \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{2n} \end{bmatrix} \quad X' = \begin{bmatrix} 1 & \dots & 1 \\ x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \end{bmatrix}$$

$$X'X = \begin{bmatrix} n & \sum x_{1i} & \sum x_{2i} \\ \sum x_{1i} & \sum x_{1i}^2 & \sum x_{1i}x_{2i} \\ \sum x_{2i} & \sum x_{1i}x_{2i} & \sum x_{2i}^2 \end{bmatrix} \quad X'Y = \begin{bmatrix} \sum y_i \\ \sum x_{1i}y_i \\ \sum x_{2i}y_i \end{bmatrix}$$

باینسی تدرک $X'X$ را می تدرک $X'Y$ فریک تدرک B برش تدرک

طریق تدرک تدرک 3×3

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \quad \text{فریک تدرک}$$

تدرک تدرک A (تدرک تدرک)

$$\det(A) = a_{11} \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} - a_{12} \begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{vmatrix} + a_{13} \begin{vmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{vmatrix}$$

$$= a_{11} [a_{22}a_{33} - a_{23}a_{32}] - a_{12} [a_{21}a_{33} - a_{23}a_{31}] + a_{13} [a_{21}a_{32} - a_{22}a_{31}]$$