

تحلیل رگرسیون

رگرسیون در واقع نوعی مدل‌سازی ریاضی است. در رگرسیون باطالع مشاهده ارتباط یک پدیده با الگو رسم که با آن الگو تناسبی داریم رفتار پدیده مورد نظر را بر داریم

$$y = f(x_1, x_2, \dots, x_n)$$

(... و تناسبی بین یک متغیر، سن، درآمد و هزینه سفر) $f =$ رضایت سفر

چون صحت از پیش بین وجود یک خطا (Error) پیش می‌آید و خطای زیاد زیانی دارد لذا باید توابع زبان را بر هم کنیم

انواع توابع زبان در مدل‌سازی

پیش‌بینی \hat{y}_i y_i : مقدار واقعی

۱- حداقل سازی مجموع قدر مطلق خطا

$$\min \sum_{i=1}^n |y_i - \hat{y}_i| = \min \sum_{i=1}^n |e_i|$$

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

$$y = \alpha x + \beta$$

رضن کنند برای یک تابع یک متغیره صورت

$$\hat{y}_i = \alpha x_i + \beta$$

$$\min \sum |y_i - \alpha x_i - \beta|$$

$$y_i - \alpha x_i - \beta = u_i - v_i \rightarrow \text{تغییر متغیره}$$

$$\min \sum_{i=1}^n u_i + v_i$$

$$y_i - \alpha x_i - \beta = u_i - v_i$$

$$u_i, v_i \geq 0$$

۱۲ الگوریتم حداقل حد اکثر الزامات

$$\min \max |e_i|$$

$$\min \max |y_i - \alpha x_i - \beta|$$

$$\max |y_i - \alpha x_i - \beta| = z$$

$$\min z$$

$$\begin{cases} y - \alpha x_i - \beta \leq z \\ y - \alpha x_i - \beta \geq -z \end{cases}$$

$$y - \alpha x_i - \beta \geq -z$$

رابطه میان سرمایه و میزان تولید برای بخش در فرمای صورت زیر می آید

	۱	۲	۳	۴	۵	۶
سرمایه	۲۵۰	۳۰۰	۴۲۰	۵۰۰	۶۰۰	۷۰۰
تولید	۵۰	۵۵	۸۰	۱۲۲	۱۵۰	۲۰۰

با استفاده از تابع زیرین مجموع قدر مطلق خطاها و ... می یابیم قدر مطلق خطاها، پارامترها مدل را برآورد کنیم

$$\min \sum_{i=1}^7 u_i + v_i + u_2 + v_2 + \dots + u_7 + v_7$$

$$250 - 250\alpha - \beta = u_1 - v_1$$

$$300 - 300\alpha - \beta = u_2 - v_2$$

}

$$700 - 700\alpha - \beta = u_7 - v_7$$

$$u_i, v_i \geq 0$$

$$\alpha = 1.25$$

$$\beta = 0$$

$$\min z$$

$$L.p$$

$$d_0 - 2d_1 \alpha - \beta \leq z$$

$$d_0 - 2d_1 \alpha - \beta \geq -z$$

$$\alpha = 0.25$$

$$\beta = 0$$

$$r_0 - 7r_1 \alpha - \rho \leq z$$

$$r_0 - 7r_1 \alpha - \rho \geq -z$$

اگر حدائق حداکثر خطا، نقطه بدترین حالت توپ می‌کنند پس آنگاه خطبه نسبت درست بدینا است

اگر مجموع قدر مطلق خطاها، درین روش خود قدر مطلق زیاد شود نسبت
اگر (در سوم) حدائق مجزوم خطا

$$\min \sum_{i=1}^m e_i^r$$

$$y = f(x_1, \dots, x_k)$$

$$= b_0 + b_1 x_1 + \dots + b_k x_k$$

$$\min \sum (y_i - b_0 - b_1 x_{i1} - b_2 x_{i2} - \dots - b_k x_{ik})^2$$

خوبنفا، اگر نوع یک اگر شفاف شد است که با هینه ساز آن می‌توانیم پارامترها
اگر ریاضی مورد نظر را برآورد کنیم شد در یک تابع یک متغیر، شکل

~~$$y = b_0 + b_1 x_1$$~~

$$b_1 = \frac{n \sum x y - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

x	0	1	1
y	2	1	2

$$\sum xy = (0 \times 2) + (1 \times 1) + (1 \times 2) = 3$$

$$\sum x^2 = 0^2 + 1^2 + 1^2 = 2$$

$$\sum x = 2$$

$$\sum y = 5$$

$$b_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} = \frac{3 \times 3 - (2)(5)}{2 - 4} = \frac{-1}{-2} = \frac{1}{2}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$b_0 = \frac{5}{3} + \frac{1}{2} \left(\frac{2}{3} \right) = \frac{5}{3} + \frac{1}{3} = \frac{6}{3} = 2$$

$$y = 2 - \frac{1}{2}x$$

کاربرد جبر ماتریس در برآورد رگرسیون حداقل مربعات خطی

$$Y = X\beta + U$$

$$X = \begin{pmatrix} 1 & x_{11} \\ \vdots & \vdots \\ 1 & x_{1n} \end{pmatrix}$$

$$\begin{pmatrix} x_{k1} \\ \vdots \\ x_{kn} \end{pmatrix}$$

$$U = \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix}$$

$$\hat{\beta} = \begin{pmatrix} b_0 \\ b_1 \\ \vdots \\ b_k \end{pmatrix}$$

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

Micro

Macro

اعتباری مدل
اعتباری مدل

۱- اعتباری مدل: عمدتاً با افراد (ضرایب) است. دانسته وجودشان در مدل ضرورت دارد یا نه؟
برای پاسخ به این سؤال از آزمون t تست برای هکت از ضرایب استفاده می شود.
مراحل اعتباری مدل:

۱- فرض ها آوری: در اعتباری مدل

$$\begin{cases} \beta_i = 0 \\ \beta_i \neq 0 \end{cases}$$
 تآید یا رد ضرایب

۲- آزمون آزمون

$$t = \frac{b_i - \beta_i}{S_{b_i}}$$

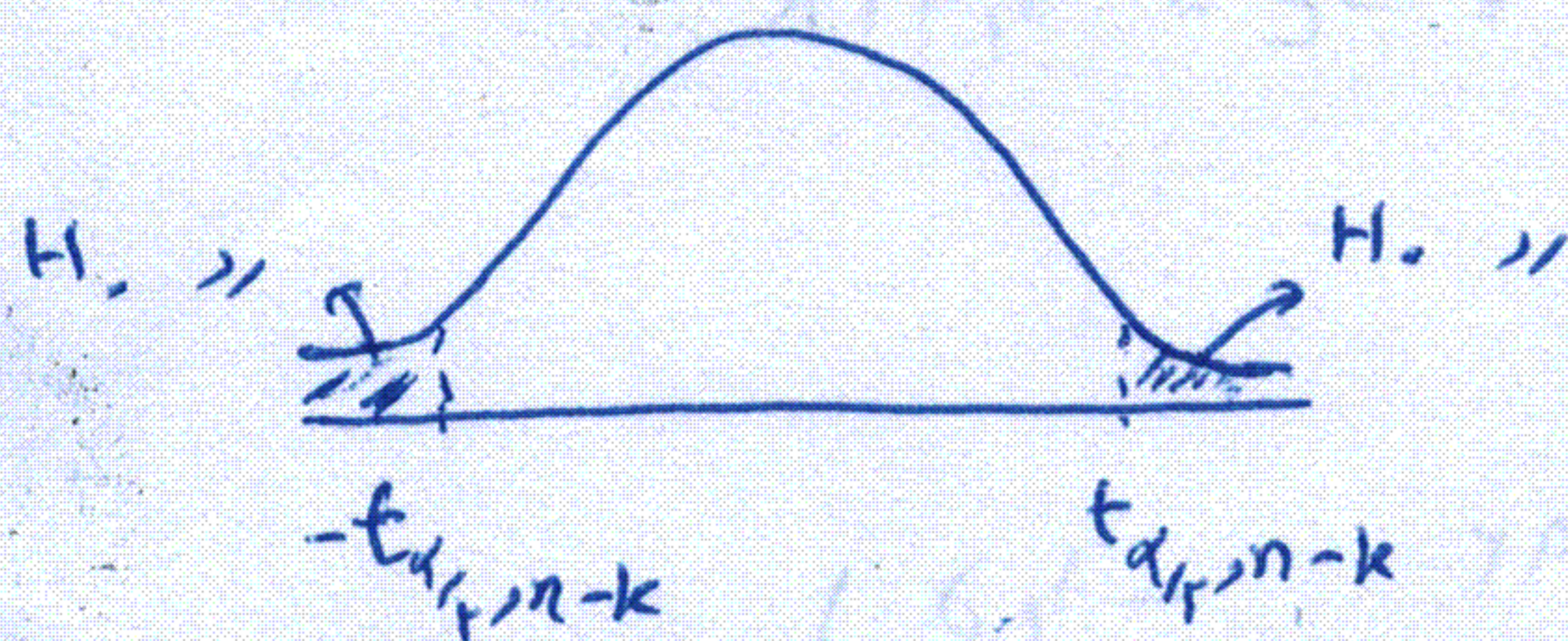
S_{b_i} - خطای معیار ضرایب

n: تعداد مشاهدات

k: تعداد ضرایب

$t_{\alpha/2, (n-k)}$

۳- مقادیر بحرانی



خطای معیار آزمون S_{b_i}

مجموع مربعات خطا

$$1) SSE = (Y^T Y) - (\hat{\beta}^T X^T Y)$$

Sum of square for error

$$SSE = \sum y^2 - b_0 \sum y - b_1 \sum xy$$

مدل

$$2) S_e^2 = \frac{SSE}{n-k}$$

S_e خطای معیار برآورد

$$3) \text{var}(\hat{\beta}) = S_e^2 (X^T X)^{-1}$$

$$\begin{pmatrix} S_{b_0}^2 & S_{b_0 b_1} & S_{b_0 b_k} \\ S_{b_1 b_0} & S_{b_1}^2 & S_{b_1 b_k} \\ S_{b_k b_0} & S_{b_k b_1} & S_{b_k}^2 \end{pmatrix}$$

حرفه بلیف	3	5	4	7	9	6	5	4	8
حرفه زارگر	11	20	16	24	26	15	21	18	27

$$(X^T X)^{-1} = \frac{1}{288} \begin{bmatrix} 321 & -51 \\ -51 & 9 \end{bmatrix}$$

$$X^T Y = \begin{pmatrix} 148 \\ 1082 \end{pmatrix}$$

$$\hat{\beta} = \begin{bmatrix} 4.79 \\ 2.292 \end{bmatrix}$$

$$\hat{y} = 4.79 + 2.292x$$

$$SSE = \sum y^2 - b_0 \sum y - b_1 \sum xy =$$

$$SSE = 3748 - (4.79 \quad 2.292) \begin{pmatrix} 148 \\ 1082 \end{pmatrix} = 59.614$$

$$(\sum y^2) - (\hat{\beta}^T X^T Y)$$

$$S_e^2 = \frac{59.614}{9-2} = 8.517$$

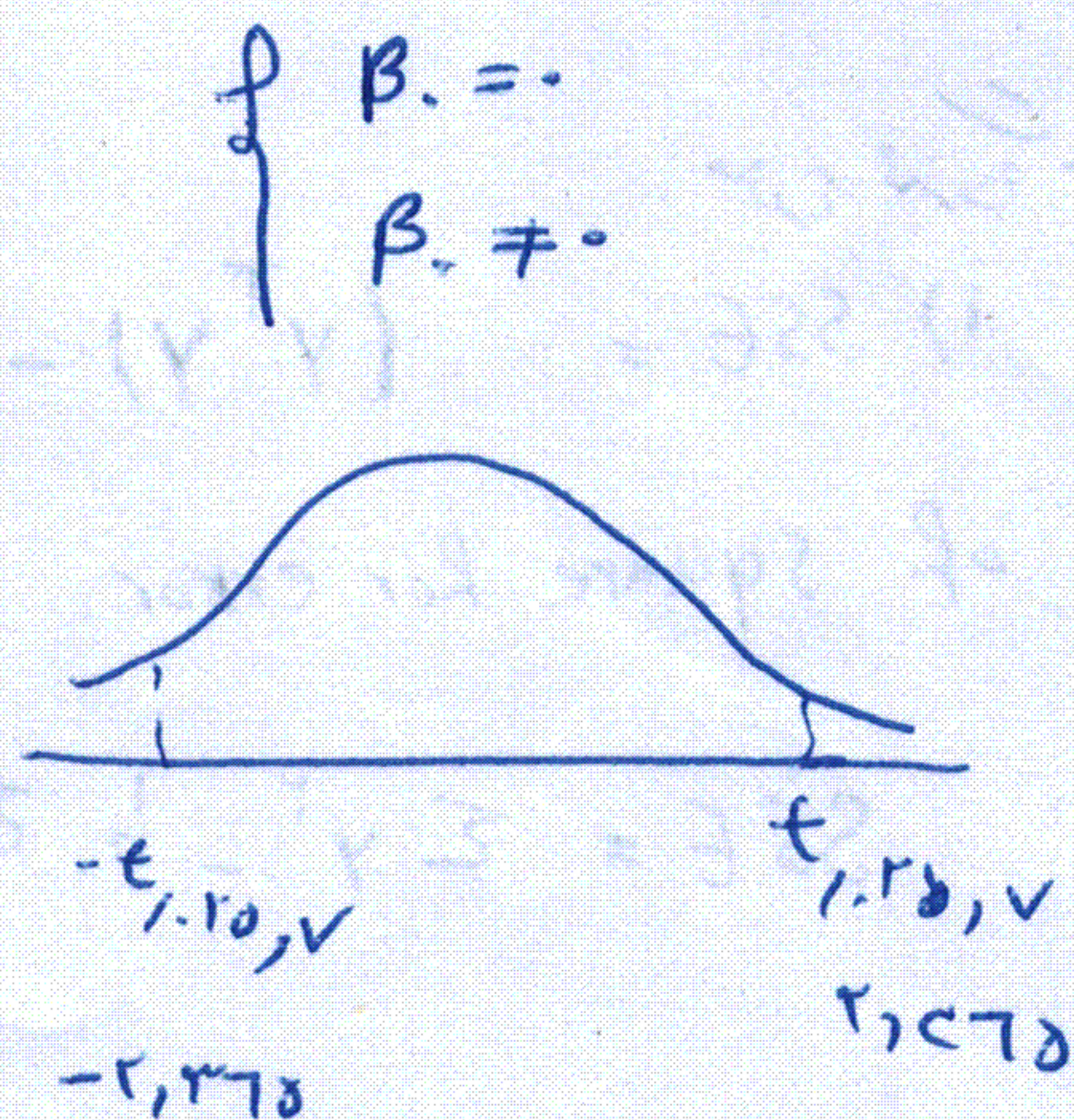
$$S_e = 2.918$$

$$\text{Var}(\hat{\beta}) = S_e^2 (X^T X)^{-1} = \begin{pmatrix} 9.892 & -1.8 \\ -1.8 & .1277 \end{pmatrix}$$

$$S_{b_0} = \sqrt{9.892}$$

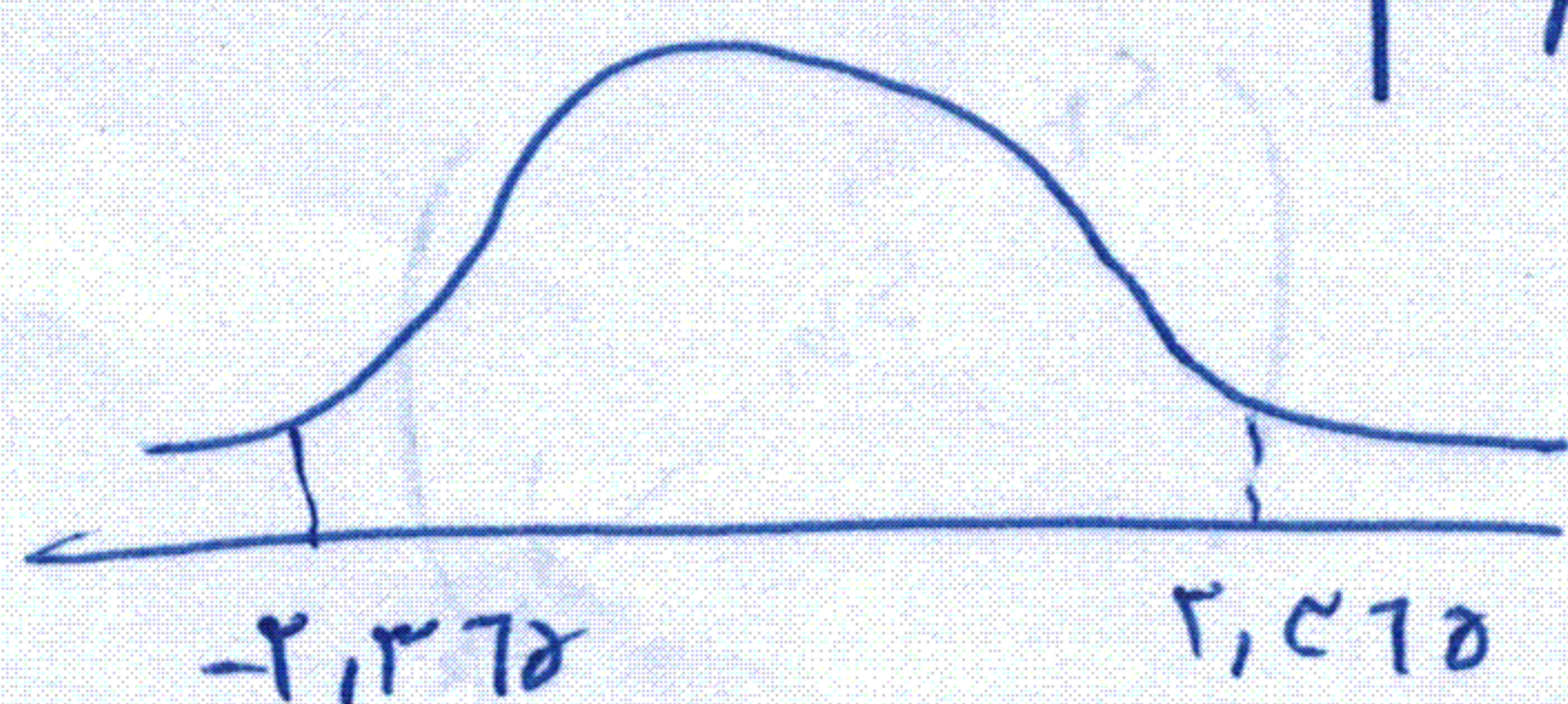
$$S_{b_1} = \sqrt{.1277}$$

$$t = \frac{4.79}{\sqrt{9.892}} = 2.25$$



$$t = \frac{2.292}{\sqrt{.1277}} = 2.052$$

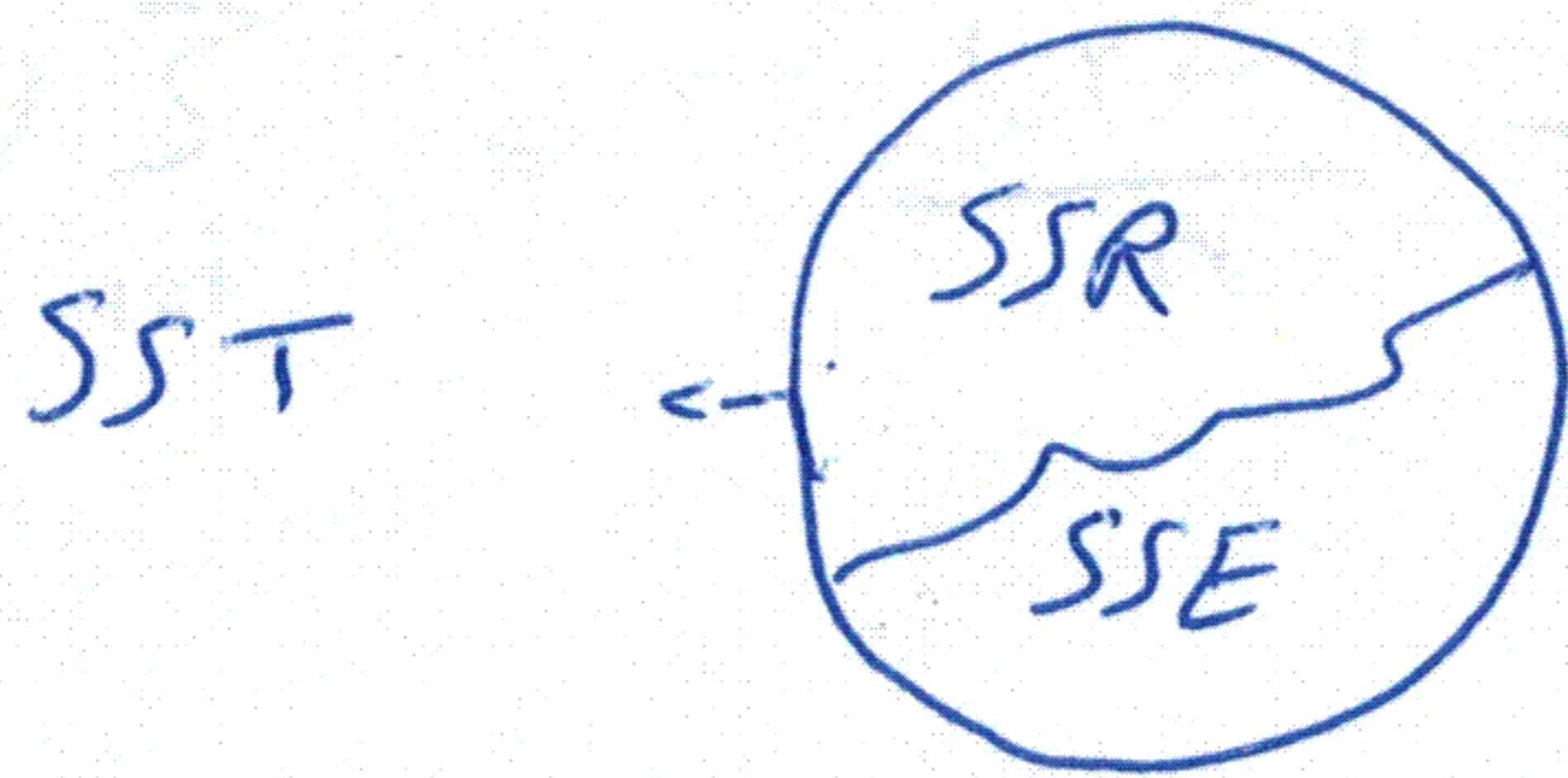
$\beta_1 = 0$
 $\beta_1 \neq 0$



H_0 و
 تأیید (غیر)

اعتبار نخبی کلان :

برای اعتبار نخبی کلان معادله برآورد شده از تکنیک تحلیل واریانس استفاده می شود. در تکنیک تحلیل واریانس، مجموع تغییرات در داده ها آماری (SST) به دو بخش مکتب افزای می شود. یکی مجموع تغییرات ناشی از مدل (SSR) و دیگری مجموع تغییرات ناشی از خطا (SSE). پس با مقایسه این دو شاخص در قالب یک آزمون آماری به نام F، معناداری کل مدل مورد بررسی قرار می گیرد.



$$SST = SSR + SSE$$

برای اعتبار نخبی کلان جدول تحلیل واریانس زیر شکل می گیرد.

منبع تغییر	مجموع مجزوات	درجه آزادی	میانگین مجزوات	F_0	F_t
رگرسیون	SSR مجموع مجزوات رگرسیون	$k-1$	$MSR = \frac{SSR}{k-1}$	$F_0 = \frac{MSR}{MSE}$	$F_{\alpha, k-1, n-k}$
خطا	SSE مجموع مجزوات خطا	$n-k$	$MSE = \frac{SSE}{n-k}$		
کل	SST مجموع مجزوات کل	$n-1$			

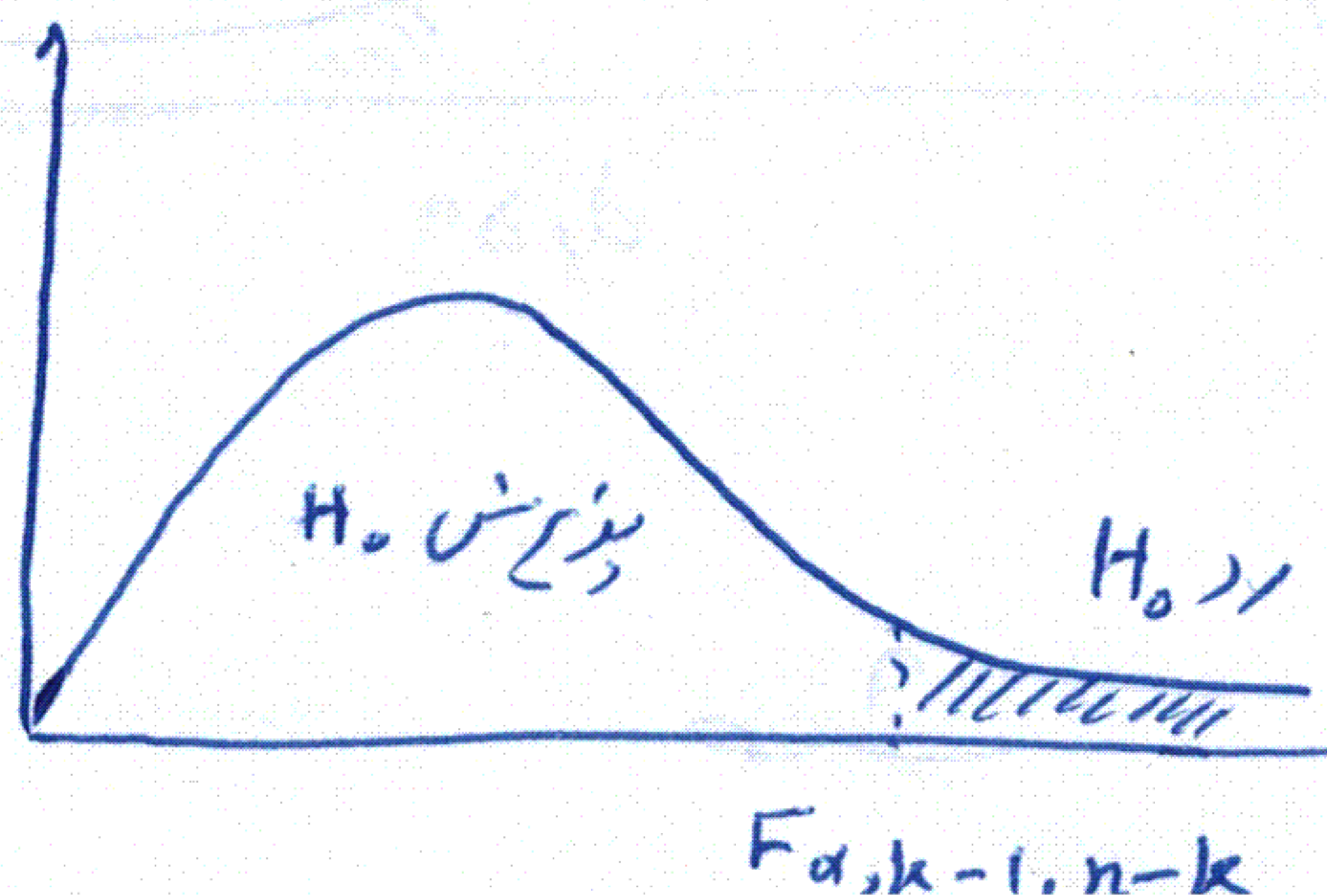
H_0 : مدل فاقد اعتبار

H_1 : مدل دارای اعتبار

$$F_0 = \frac{MSR}{MSE}$$

در H_0 : اعتبار مدل $F_0 > F_t$

در H_1 : اعتبار مدل $F_0 \leq F_t$



X	۳	۵	۴	۷	۹	۶	۵	۴	۸
Y	۱۱	۲۰	۱۷	۲۲	۲۷	۱۵	۲۱	۱۸	۲۷

$$SST = \sum y^2 - \frac{(\sum y)^2}{n} = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$SSE = \sum y^2 - b \cdot \sum y - b_1 \sum xy = \sum_{i=1}^n (\hat{y}_i - \hat{y})^2$$

$$SSR = SST - SSE = b \cdot \sum y + b_1 \sum xy - \frac{(\sum y)^2}{n} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

در صورت

$$SST = ۳۷۴۸ - \frac{1}{9} (14۸)^2 = ۲۲۷,۵۵۵$$

$$SSE = ۳۷۴۸ - ۶,۷۸۹ (14۸) - ۲,۲۹۲ (1.۸۲) = ۵۹,۷۱۲$$

$$SSR = SST - SSE = ۱۶۷,۸۴۳$$

$$df_1 = k - 1 = ۲ - 1 = 1$$

$$df_2 = n - k = ۹ - ۲ = ۷$$

$$MSR = \frac{SSR}{k-1} = \frac{167,843}{1} = 167,843$$

$$MSE = \frac{SSE}{n-k} = \frac{59,712}{7} = 8,530$$

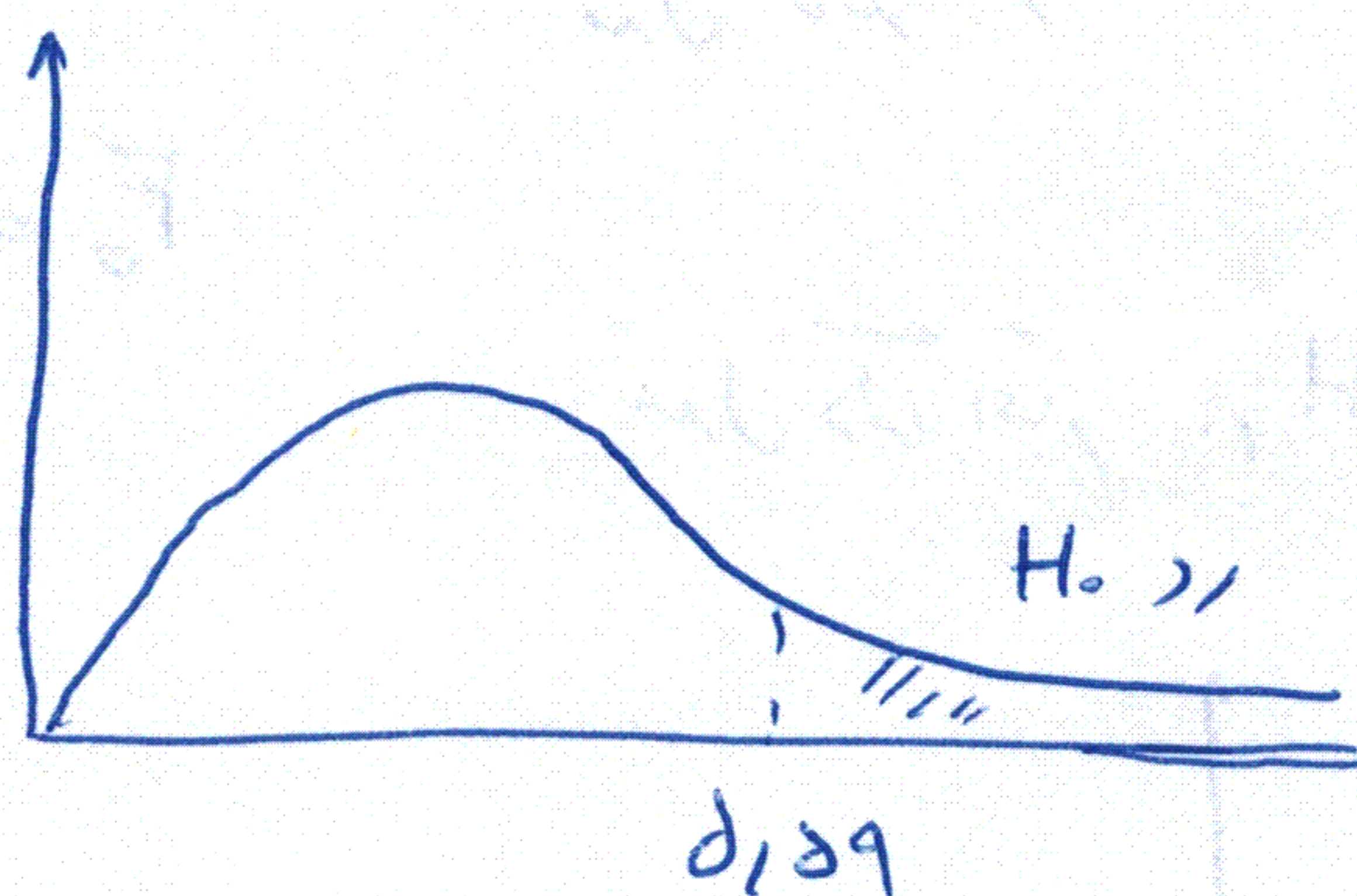
$$F_0 = \frac{MSR}{MSE} = 19,52$$

$$F_{\alpha} = F_{0,05,1,7} = 5,59$$

$$F_0 > F_{\alpha}$$

H₀ رد

تفاوت معنی دار



ضریب تعیین (R^2)

این ضریب نشان دهنده میزان قدرت تبیین کنندگی مدل است. نشان می‌دهد تا چه حد تغییرات متغیر وابسته توسط متغیرها مستقل قابل تغییر و تغییر است

$$(1) \quad SST = SSE + SSR$$

(2) فرض بر SST در تقسیم شوند $\rightarrow \frac{SST}{SST} = \frac{SSE}{SST} + \frac{SSR}{SST} \rightarrow R^2$

$$(3) \quad 1 = \frac{SSE}{SST} + R^2$$

$$(4) \quad R^2 = 1 - \frac{SSE}{SST} \quad \text{و} \quad R^2 = \frac{SSR}{SST}$$

مثال قبل $R^2 = \frac{178}{227.5} = 0.78$

یعنی تقریباً ۷۸٪ از تغییرات فروش ناشی از تغییرات تبلیغات است و ۲۲٪ بقیه ناشی از عوامل ناشناخته است که هنوز مطالعه نشده ایم.