

1) Testes não paramétricos:

→ Teste de Aderência: Queremos estimar o tipo de distribuição que a população segue a partir dos dados de uma amostra (dados em "box" ou "na" aderência ao modelo)
 A hipótese testada refere-se a forma de distribuição da população.

$$\begin{cases} H_0: \text{os dados observados provêm de uma distribuição } "X" \text{ com parâmetro } "Y" \\ H_1: \text{os dados observados NÃO provêm de uma distribuição } "X" \text{ com parâmetro } "Y" \end{cases}$$

• Teste de Aderência pelo χ^2

$$\chi^2_{\text{calc}} = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} = \sum_{i=1}^k \frac{O_i^2}{E_i} - n$$

→ média dos O_i calculado a partir da distribuição que estamos analisando com o dado parâmetro ($E_i = n \cdot P_i(X=x_i)$)
 poisso, bimodal, etc.

Obs: $E_i \geq 5$, caso contrário somar os valores de $E_i < 5$ com a classe anterior em que $E_i \geq 5$.
 (fazer o mesmo com as observações)

→ observações realizadas = frequências

$$g.l. = r - (k - 1) - m$$

→ número de parcelas após a junção

→ número de parâmetros estimados a partir dos dados da amostra (para o cálculo dos χ^2).

$$n = \sum_{i=1}^k E_i = \sum_{i=1}^k O_i$$

Se $\chi^2_{\text{calc}} > \chi^2_{r, \alpha} \Rightarrow$ rejeito H_0 !

• Teste de Aderência pelo "Papel de Probabilidade Normal"

→ para estimar μ tomamos o ponto de 50%
 → para estimar σ tomamos -1 σ e +1 σ

No eixo horizontal devemos colocar as frequências (observações) e no eixo vertical devemos colocar as frequências relativas acumuladas dada por:

$$N = \frac{50(2i-1)}{n} \quad i = 1, \dots, n$$

→ Se obtivermos uma reta, os dados tem aderência pela normal!

• Teste de Aderência pelo Método de Kolmogorov - Smirnov

→ ordenar os dados da amostra!

Trata-se de um Teste mais poderoso que o χ^2 para testar a aderência. A variável de teste é a maior diferença observada entre a função de distribuição acumulada do modelo e da amostra.

→ f.d.a do modelo: $F(x) = P(X \leq x)$ → utilizamos a normal!

→ f.d.a da amostra: $G(x)_{\text{esquerda}} = \frac{(i-1)}{n}$ $Z_i = \frac{x_i - \mu}{\sigma}$
 $G(x)_{\text{direita}} = \frac{i}{n}$

$$d_{\text{calc}} = \max |F(x) - G(x)|$$

Obs: se $n > 50$: $d_{\text{crit}} = \frac{1,36}{\sqrt{n}}$ e $d_{\text{crit}} = \frac{1,63}{\sqrt{n}}$
 $\alpha = 5\%$ $\alpha = 1\%$

→ "se for pequena suficiente" → "rejeito" H_0

Se $d_{\text{calc}} > d_{\text{crit}} \Rightarrow$ rejeito H_0 !



→ Teste de Independência:

- H_0 : as variáveis são independentes
- H_1 : as variáveis não são independentes

$$\chi^2_v = \chi^2_{calc} = \sum_{i=1}^r \sum_{j=1}^k \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$E_{ij} = \frac{j_i \cdot b_j}{n}$$

$$v = (r-1)(k-1)$$

Inquiência observada na interseção da linha i com a coluna j

II) ANOVA

→ Análise de Variância (ANOVA)

→ comparar as médias das populações ⇒ identifica diferenças entre médias populacionais devido a diversas causas (classificações) atuando simultaneamente sobre os elementos da população.

- Hipóteses
- Populações homocedásticas
 - Populações com distribuições iguais normais

→ variâncias iguais

1) Uma classificação - amostras de mesmo tamanho

- $H_0: \mu_1 = \mu_2 = \dots = \mu_k$
- H_1 : pelo menos uma das médias é diferente das demais

Notação:

$$T_i = \sum_{j=1}^n x_{ij} \rightarrow T = \sum_{j=1}^k T_j$$

$$Q_i = \sum_{j=1}^n x_{ij}^2 \rightarrow Q = \sum_{j=1}^k Q_j$$

Posso subtrair o menor termo para facilitar os cálculos!
 $\sigma^2(c + x) = \sigma^2(x)$

k = número de amostras
 n = tamanho da amostra
 $N = n \cdot k$ = tamanho total

Fonte de Variação	Soma de Quadrados	G.L.	Quadrado Médio	F _{calc}	F _α
Entre Amostras	$SQE = \frac{\sum_{i=1}^k T_i^2}{n} - \frac{T^2}{nk}$	$k-1$	$S_E^2 = \frac{SQE}{k-1}$	$F = \frac{S_E^2}{S_R^2}$	$F_{(k-1, k(n-1), \alpha)}$
Residual	$SQR = Q - \frac{\sum_{i=1}^k T_i^2}{n}$	$k(n-1)$	$S_R^2 = \frac{SQR}{k(n-1)}$		
Total	$SQT = Q - \frac{T^2}{nk}$	$nk-1$			

→ numerador
 → denominador

→ Sendo H_0 verdadeiro podemos estimar σ^2 com um de 3 maneiras possíveis: S_T^2, S_E^2, S_R^2

- $H_0: \sigma_E^2, \sigma_R^2$
- $H_1: \sigma_E^2 \Rightarrow \sigma_R^2$

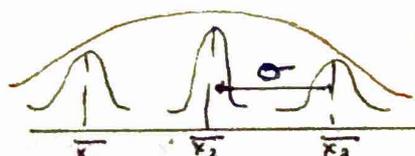
Espace Fundamental: $SQT = SQR + SQE$

Probabilidades que não sabemos aplicar/separar (entre elementos de cada amostra ⇒ dentro da amostra)

Para amostras de mesmo tamanho: $S_R^2 = \frac{S_1^2 + S_2^2 + \dots + S_k^2}{k}$

Se não há variações dentro das amostras e entre as amostras podemos estimar que as amostras provêm da mesma população.

Se $F_{calc} > F_{crit} \Rightarrow$ rejeito H_0



→ H_0 é verdadeiro

2) Uma classi ficao - amostras de tamanhos diferentes

Fonte de Variao	Soma de Quadrados	G.L.	Quadrado Mdio	Fonte	F _α
Entre Amostras	$SQE = \sum_{i=1}^k \frac{T_i^2}{n_i} - \frac{T^2}{\sum_{i=1}^k n_i}$	k-1	$S_E^2 = \frac{SQE}{k-1}$	F _{α, k-1, \sum n_i - k}	
Residual	$SQR = Q - \sum_{i=1}^k \frac{T_i^2}{n_i}$	$\sum_{i=1}^k n_i - k$	$S_R^2 = \frac{SQR}{\sum n_i - k}$		
Total	$SQT = Q - \frac{T^2}{\sum_{i=1}^k n_i}$	$\sum_{i=1}^k n_i - 1$			

$$SQT = \sum x^2 - \frac{G^2}{L}$$

$$S_R^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2 + \dots + (n_k-1)s_k^2}{n_1 + n_2 + \dots + n_k - k}$$

Quais mdias devem ser consideradas diferentes de quais outras?

Mtodo de Tukey: para k linhas e n colunas, as mdias \bar{x}_i e \bar{x}_m srio consideradas distintas se:

$$|\bar{x}_i - \bar{x}_m| > t_{k, \alpha} \sqrt{\frac{S_R^2}{n}}$$

para cada amostra

amplitude studentizada

Mtodo de Scheffé: as mdias \bar{x}_i e \bar{x}_m srio consideradas distintas entre si, se:

$$|\bar{x}_i - \bar{x}_m| > \sqrt{S_R^2 \cdot \frac{2(k-1)}{n} F_{k-1, k(n-1), \alpha}}$$

Dois Classi ficao: os elementos observados srio classificados segundo dois critrios. Admitindo que haja nk observaes, do primeiro critrio temos n elementos e k amostras, e do segundo critrio temos k elementos e duas amostras.

$$H_{01}: \mu_{1,linha} = \mu_{2,linha} = \dots = \mu_{k,linha}$$

$$H_{02}: \mu_{1,coluna} = \mu_{2,coluna} = \dots = \mu_{n,coluna}$$

H_1 : pelo menos uma das médias é diferente das demais \rightarrow linha interfere

H_2 : pelo menos uma das médias é diferente das demais \rightarrow coluna interfere

Modelo Fixo: os efeitos resultantes das classificações segundo linhas e colunas são ambos fixos, ou seja representam a totalidade das condições existentes

3) Duas Classificações sem repetição:

\rightarrow Modelo fixo \Rightarrow inexistência de interação entre linhas e colunas.

\rightarrow podemos estimar σ^2 através de S_T^2, S_L^2, S_C^2 e S_R^2 .

Fonte de Variação	Soma de Quadrados	Gr. L.	Quadrado Médio	F_{calc}	F_{c}
Entre Linhas	$S_{QL} = \frac{\sum T_i^2}{n} - \frac{T^2}{nk}$	$k-1$	$S_L^2 = \frac{S_{QL}}{k-1}$	$F_L = \frac{S_L^2}{S_R^2}$	$F_{k-1, (k-1)(n-1), \alpha}$
Entre Colunas	$S_{QC} = \frac{\sum T_j^2}{k} - \frac{T^2}{nk}$	$n-1$	$S_C^2 = \frac{S_{QC}}{n-1}$	$F_C = \frac{S_C^2}{S_R^2}$	$F_{n-1, (k-1)(n-1), \alpha}$
Residual	$S_{QR} = S_{QT} - S_{QC} - S_{QL}$	$(k-1)(n-1)$	$S_R^2 = \frac{S_{QR}}{(k-1)(n-1)}$		
Total	$S_{QT} = Q - \frac{T^2}{nk}$	$nk-1$			

$k =$ linhas

$n =$ colunas

4) Duas Classificações com repetição:

\rightarrow se houver interação (testa o efeito fixo)

Fonte de Variação	Soma de Quadrados	Grupos de Liberdade	Quadrado Médio	F_{calc}	F_c
Entre Linhas	$S_{QL} = \frac{\sum T_i^2}{nr} - \frac{T^2}{nkr}$	$k-1$	$S_L^2 = \frac{S_{QL}}{k-1}$	$F_L = \frac{S_L^2}{S_R^2}$	
Entre Colunas	$S_{QC} = \frac{\sum T_j^2}{kr} - \frac{T^2}{nkr}$	$n-1$	$S_C^2 = \frac{S_{QC}}{n-1}$	$F_C = \frac{S_C^2}{S_R^2}$	
Interação	$S_{QS} = S_{QT} - S_{QL} - S_{QC}$	$(k-1)(n-1)$	$S_S^2 = \frac{S_{QS}}{(k-1)(n-1)}$	$F_S = \frac{S_S^2}{S_R^2}$	$F_{(k-1)(n-1), nk(r-1), \alpha}$
Entre Tratamentos	$S_{QTr} = \frac{\sum \sum T_{ij}^2}{r} - \frac{T^2}{nkr}$	$nk-1$	$S_{Tr}^2 = \frac{S_{QTr}}{nk-1}$	$F_{Tr} = \frac{S_{Tr}^2}{S_R^2}$	$F_{nk-1, nk(r-1), \alpha}$
Residual	$S_{QR} = S_{QT} - S_{QTr}$	$nk(r-1)$	$S_R^2 = \frac{S_{QR}}{nk(r-1)}$		
Total	$S_{QT} = Q - \frac{T^2}{nkr}$	$nkr-1$			

"Tratamento" é cada combinação entre linha e coluna.

↳ pode haver uma variação entre os nk tratamentos aos quais cada elemento é submetido

Portanto, também podemos estimar σ^2 por $\hat{\sigma}_r^2$.

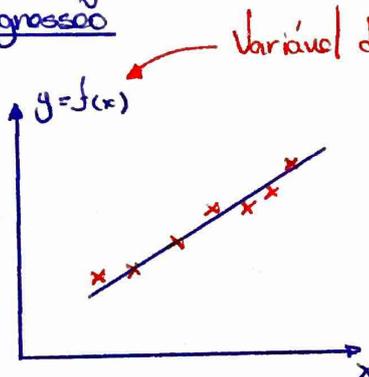
"Interação" está ligada ao fato que os critérios se relacionam

↳ Devemos verificar se há ou não interação entre os critérios (linhas e colunas). Se não

houver, o resíduo "absorve a interação" (somar os SQR e os G.L.)

Depois recalcular $\hat{\sigma}_r^2$

III) Correlação e Regressão



Variável dependente

Regressão Linear: trata-se da reta que minimiza a distância dos pontos em y.

Aplicar o conceito de estratificação (dados provenientes de conjuntos homogêneos)

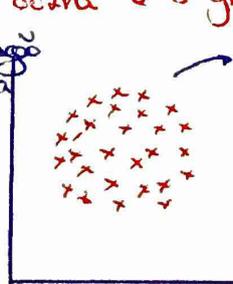
Variável independente (temos total controle sobre ela)

Correlação Linear:

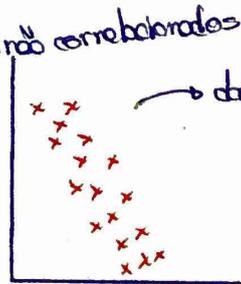
↳ nos mostra como as variáveis tendem a se relacionar fornecendo a magnitude do crescimento que uma possui em relação a outra e o grau de relacionamento



$-1 < r < 0$



$r \approx 0$



$0 < r < 1$

dados com correlação negativa

Obs:

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

Coefficiente de Pearson:

é sobre p/rela!

$$r = \frac{\sum xy}{\sqrt{\sum x^2 \cdot \sum y^2}}$$

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

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

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

$-1 \leq r \leq +1$

se $r=1$ ou $r=-1$ temos que as variáveis se relacionam por uma reta

→ pode ser ⊕ ou ⊖

→ ⊕ sempre

→ ⊕ sempre

Coefficiente R^2 :

$$R^2 = \frac{SQT - \sum QR}{SQT} = r^2$$

é mais geral

Coeff. de Correlação populacional

Teste para o coeficiente de correlação:

$H_0: \rho = 0$ $t_{calc} = r \sqrt{\frac{n-2}{1-r^2}}$

$H_1: \rho \neq 0$ $t_{crit} = t_{n-2, \alpha/2}$

→ espaço que relaciona as variáveis de interesse.

→ Regressão → após verificada a correlação entre os dados, faz-se necessário encontrar uma função que exprime

esse não é o caso

dimos o total controle sobre as

Hipóteses: x_1, \dots, x_n são admitidos sem erros

y é admitido com erro. Este erro é $\epsilon \sim N(0, \sigma^2)$

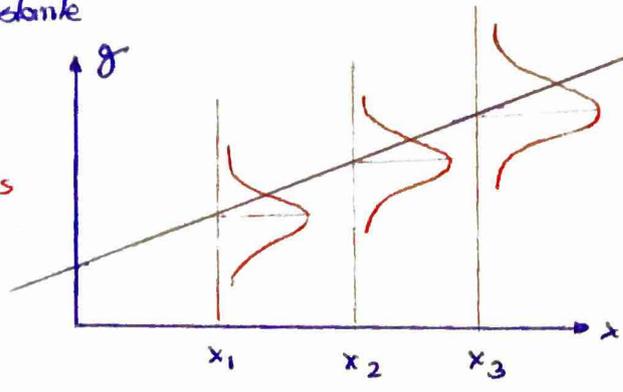
$\sigma^2 = \text{constante}$

→ média 0

→ variância residual

Regressão Linear Simples:

Quas variáveis



Modelo Teórico:

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

média de y para cada x

Rela Teórica: $y_i = \alpha + \beta x$ → rela verdadeira

Rela Estimativa: $\hat{y} = a + bx$

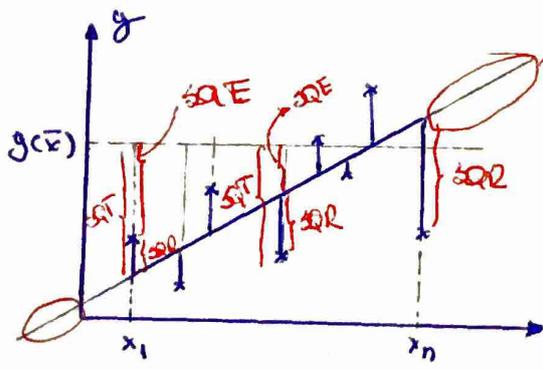
coef. linear

coef. angular

Para achar a e b utilizaremos o método dos mínimos quadrados (MMQ)

Método dos Mínimos Quadrados:

a e b que mais se ajustam ao modelo são aqueles que tornam mínimo a diferença entre cada ponto e a rela estimada (distância mínima).



→ não sabemos o que ocorre

$$SQTy = \sum (y_i - \bar{y})^2$$

$$SQR_{g}^{\min} = \min \sum d_i^2 =$$

$$= \min \sum (y_i - \hat{y}_i)^2 =$$

$$= \min \sum (y_i - a - bx_i)^2$$

o MMQ vai reduzir ao máximo a variância residual (SQR)

$$b = \frac{\sum xy}{\sum xx}$$

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

Obs:

- Sempre há uma rela que representa um conjunto de dados
- Mas nem sempre esta rela é o melhor modelo

$$\hat{y} = a + bx$$

Valor Médio!

Variância Residual:

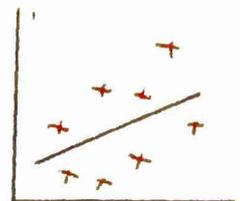
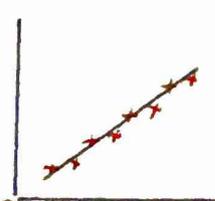
$$s_e^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2} = \frac{\sum yy - b \sum xy}{n-2}$$

→ estimador de a e b

Variância em torno da reta dos mínimos quadrados.

$$SQR + \sum (\hat{y}_i - \bar{y})^2 = \sum yy - SQT$$

$$\sum (y_i - \hat{y}_i)^2$$



Variância Residual Pequena

→ s_e^2 grande

a e b são variáveis aleatórias $\rightarrow a, b \sim \text{Normal}(\mu, \sigma)$

Variável b :

$$\mu(b) = \beta$$

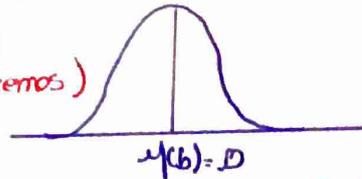
$$\sigma^2(b) = \frac{\sigma_e^2}{\sum x_i^2}$$

$$\hat{\sigma}^2(b) = \frac{\hat{\sigma}_e^2}{\sum x_i^2}$$

Escremos supondo que o modelo linear é satisfatório

Resíduos (não conhecemos)

Maiores amostra, menor será o $\sigma^2(b)$



Intervalos de confiança

$a, b, \hat{y} \rightarrow \alpha, \beta, \hat{y}$
média \rightarrow previsões

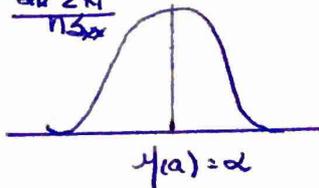
Variável a :

$$\mu(a) = \alpha$$

$$\sigma^2(a) = \frac{\sigma_e^2 \sum x_i^2}{n \sum x_i^2}$$

$$\hat{\sigma}^2(a) = \frac{\hat{\sigma}_e^2 \sum x_i^2}{n \sum x_i^2}$$

Maiores amostra, menor será o $\sigma^2(a)$



quanto mais longe x_i maior o $\hat{\sigma}^2$!

Intervalo de Confiança (para $y_i = \alpha + \beta x_i$)

$$\mu(y_i) = \alpha + \beta x_i$$

$$\text{Estimativa: } \alpha + \beta x_i$$

onde a relação é verdadeira dado x_i

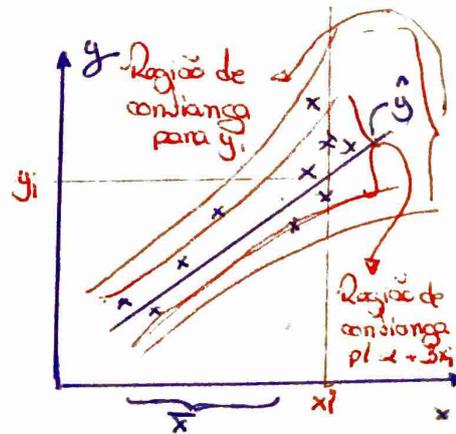
$$\text{Var}(\alpha + \beta x_i) = \sigma_e^2 \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum x_i^2} \right)$$

$$\text{IC: } \hat{y}_i \pm t_{n-2; \alpha/2} \sqrt{\hat{\sigma}_e^2 \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum x_i^2} \right)}$$

Região de confiança para $\alpha + \beta x_i$ (relação verdadeira)

ponto qualquer!

quanto mais afastado da média dos "dados controlados" maior será a variância da relação estimada!



Intervalo de Confiança (para uma previsão de y_i dado x_i)

$$\text{Estimativa: } \alpha + \beta x_i \quad \text{Var}(\alpha + \beta x_i + \epsilon_i) = \sigma_e^2 \left(1 + \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum x_i^2} \right)$$

$$\text{IC: } \hat{y}_i \pm t_{n-2; \alpha/2} \sqrt{\hat{\sigma}_e^2 \left(1 + \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum x_i^2} \right)}$$

onde y_i é observado dado x_i

Teste de Hipótese e Análise de Variância aplicadas à regressão:

Quando desejamos verificar se determinada relação de regressão tem inclinação β_0 podemos fazer o seguinte teste:

$$\begin{cases} H_0: \beta = \beta_0 \\ H_1: \beta > \beta_0 \end{cases}$$

$$t_{\text{calc}} = \frac{b - \beta_0}{\sqrt{\hat{\sigma}^2(b)}} = \frac{b - \beta_0}{\hat{\sigma}_e / \sqrt{\sum x_i^2}}$$

$$t_{\text{crit}} = t_{n-2; \alpha \text{ ou } \alpha/2}$$

teste de correlação ($\beta_0 = 0$)

Caso particular (teste de correlação):

$$\begin{cases} H_0: \beta = 0 \rightarrow \text{se } H_0 \text{ for verdadeiro então } y \text{ não depende de } x, \text{ pois, } y = \alpha + \beta x + \epsilon_i \\ H_1: \beta \neq 0 \rightarrow \text{Modelo tem validade (há correlação)} \end{cases}$$

há correlação

Quando queremos saber se a reta passa pela origem : $H_0: \alpha = 0$ (não passa pela origem)
 $H_1: \alpha \neq 0$

$t_{calc} = \frac{a - \alpha_0}{\hat{\Delta}(a)}$, onde $\hat{\Delta}(a) = \frac{\hat{\Delta}_R^2 \sum x_i^2}{n \hat{\Delta}_{xx}}$ $t_{crit} = |t_{n-2, \alpha/2}|$

ANOVA para regressão:

Lembrar que: $\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$
 (semelhante a dizer o teste $B_0 = 0$ proibicionada)
 SQT (variância em torno da média)
 SQR (proporcionará a variância em torno da reta de regressão)
 SRE (proporcionará a variância explicada pela regressão)
 Verificar se o modelo de regressão é satisfatório
 Indução: se a variância em torno da média (S_y^2) for muito maior que a variância em torno da reta de regressão (S_e^2) \Rightarrow a regressão é válida.

Fórmulas:

$$SQT = \sum (y_i - \bar{y})^2 = S_{yy}$$

$$SRE = \sum (\hat{y}_i - \bar{y})^2 = b^2 \Delta_{xx} = b \Delta_{xy}$$

$$SQR = \sum (\hat{y}_i - y_i)^2 = SQT - SRE = S_{yy} - b^2 \Delta_{xx}$$

$F_{calc} > F_{crit}$: há evidências para validar a regressão.

Fonte Variância	Soma de Quadrados	G.L.	Quadrado Médio	F _{calc}	F _{crit}
Explicada pela regressão	SRE	1	$\frac{b^2 \Delta_{xx}}{1}$	$F_{calc} = \frac{b^2 \Delta_{xx}}{\Delta_R^2}$	$F_{1, n-2, \alpha}$
Residual	SQR	n-2	Δ_R^2		
Total	SQT	n-1			

Se $\sigma_E \gg \sigma_R$ a variância explicada pela regressão não provém da variância residual. Logo a uma outra razão para sua existência. Esta razão é a dependência de x na forma linear.

Regressão Polinomial

Trabalha-se de um polinômio de segundo grau: $y = a + bx + cx^2$

Estimativa: $\hat{y}_0 = a + bx + cx^2$

Para achar a, b, c basta resolver:

$$\hat{y}_0 = a + b(x - \bar{x}) + c(x - \bar{x})^2$$

$$\begin{cases} \sum y_i = na + c \sum (x_i - \bar{x})^2 - S_{xx} \\ \sum (x_i - \bar{x}) y_i = b \sum (x_i - \bar{x})^2 + c \sum (x_i - \bar{x})^3 \\ \sum (x_i - \bar{x})^2 y_i = a \sum (x_i - \bar{x})^2 + c \sum (x_i - \bar{x})^4 \end{cases}$$

Podem ser calculado igual na reta

Trabalha-se de um polinômio de grau maior que 2:

$$\begin{cases} \sum y_i = na + b \sum x_i + c \sum x_i^2 + d \sum x_i^3 + \dots \\ \sum x_i y_i = a \sum x_i + b \sum x_i^2 + c \sum x_i^3 + d \sum x_i^4 + \dots \\ \sum x_i^2 y_i = a \sum x_i^2 + b \sum x_i^3 + c \sum x_i^4 + d \sum x_i^5 + \dots \\ \vdots \\ \vdots \\ \vdots \end{cases}$$

→ regressão múltipla

↳ y agora depende de mais de um x

$$\hat{y} = a + b_1 x_1 + b_2 x_2 + \dots + b_k x_k$$

$$\begin{cases} \hat{s}_{1y} = b_1 s_{11} + b_2 s_{12} + \dots + b_k s_{1k} \\ \hat{s}_{2y} = b_1 s_{21} + b_2 s_{22} + \dots + b_k s_{2k} \\ \vdots \\ \hat{s}_{ky} = b_1 s_{k1} + b_2 s_{k2} + \dots + b_k s_{kk} \end{cases}$$

↳ calcular como s_{xx}

↳ calcular como s_{xy}

$$s_{kj} = \sum_{j=1}^n (x_{kj} - \bar{x}_k)(y_j - \bar{y}) = \sum x_{kj} y_j - \frac{\sum x_{kj} \sum y_j}{n}$$

$$\text{e.k.: } s_{kk} = \sum x_{kj}^2$$