

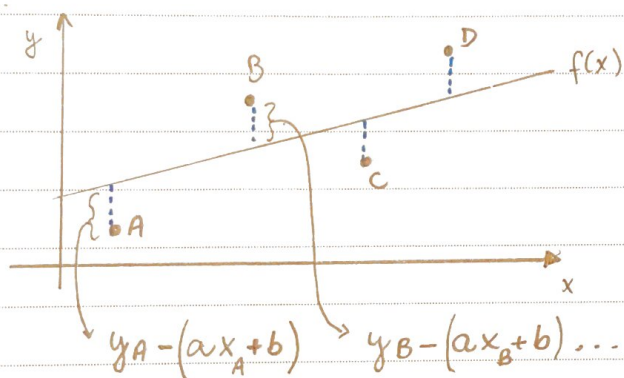
## REGRESSÃO LINEAR

• FREDERICO VILELA CURTI e RAPHAEL COSTA

I Podemos encontrar a função que melhor descreve uma distribuição de uma variável aleatória, para uma reta, generalizando-a na forma

$$f(x) = ax + b$$

Os coeficientes  $a$  e  $b$  podem ser obtidos com as seguintes etapas:



I. Encontramos o valor do erro (diferença entre o valor de  $y$  e  $f(x)$ ) em cada ponto

Como alguns erros serão negativos, elevamos todos os erros ao quadrado para contarmos isso, obtendo que o erro total acumulado assume:

$$\rightarrow E^2 = (y_A - (ax_A + b))^2 + (y_B - (ax_B + b))^2 + \dots + (y_n - (ax_n + b))^2$$

\* Simplificando a expressão, através dos produtos notáveis:

$$E^2 = y_A^2 - 2y_A(ax_A + b) + (ax_A + b)^2 + y_B^2 - 2y_B(ax_B + b) + (ax_B + b)^2 + \dots + y_n^2 - 2y_n(ax_n + b) + (ax_n + b)^2$$

\* Repletando o processo e aplicando propriedades distributivas

$$= y_A^2 - 2y_Aax_A - 2y_Ab + a^2x_A^2 + 2ax_Ab + b^2 + \dots + y_n^2 - 2y_nax_n - 2y_nb + a^2x_n^2 + 2ax_nb + b^2$$

II. Colocamos nossos termos de interesse em evidência:

$$E^2 = \left[ \left( \underline{y_A^2 + y_B^2 + \dots + y_n^2} \right) - 2a \cdot (X_A y_A + X_B y_B + \dots + X_n y_n) - \right. \\ \left. 2b (y_A + y_B + \dots + y_n) + a^2 \cdot (X_A^2 + X_B^2 + \dots + X_n^2) \right. \\ \left. + 2 \cdot ab \cdot (X_A + X_B + \dots + X_n) + n \cdot b^2 \right] \quad (I)$$

→ Os termos de  $x$  e  $y$  da expressão podem ser simplificados para suas médias  
 por exemplo:

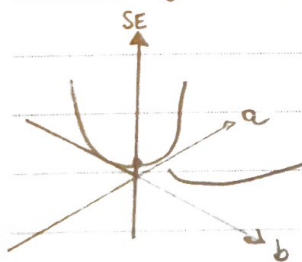
$$* \frac{y_A^2 + y_B^2 + \dots + y_n^2}{n} = \overline{y^2} \cdot n$$

II. Obtemos a fórmula:

$$\rightarrow (n \cdot \overline{y^2}) - 2an \overline{Xy} - 2bn \cdot \overline{y} + a^2 \cdot n \overline{x} + 2abn \overline{x} + nb^2$$

Cl ena expressão, temos todos os valores positivos para  $n$  e  $b$ . Com as derivadas parciais, podemos encontrar o menor valor destes coeficientes, minimizando o erro

$$E^2 = SE \text{ (Squared Error)}$$



$$\frac{\partial SE}{\partial a} = \frac{\partial SE}{\partial b} = 0$$

IV Derivada parcial

$$\frac{\partial SE}{\partial a} = -2n \overline{Xy} + 2an \overline{x^2} + 2bn \overline{x} = 0 \quad (\div 2n) \\ \rightarrow -\overline{Xy} + a \overline{x^2} + b \overline{x} = 0 \rightarrow \boxed{a \overline{x^2} + b \overline{x} = \overline{Xy}} \quad (\div \overline{x}) \rightarrow$$

$$\frac{\partial SE}{\partial b} = -2n \overline{y} + 2an \overline{x} + 2nb = 0 \quad (\div 2n) \\ \rightarrow -\overline{y} + a \overline{x} + b = 0 \rightarrow \boxed{a \overline{x} + b = \overline{y}} \rightarrow$$

$$\hat{y} = ax + b \quad \rightarrow \quad \left( \frac{\bar{x}, \bar{xy}}{\bar{x}} \right)$$

$$\quad \quad \quad \left( \bar{x}, \bar{y} \right)$$

V. Simplificando para isolar a em função de x e y

$$a \cdot \left( \bar{x} - \frac{\bar{x}^2}{\bar{x}} \right) = \bar{y} - \frac{\bar{xy}}{\bar{x}} \rightarrow a = \frac{\bar{y} - \frac{\bar{xy}}{\bar{x}}}{\bar{x} - \frac{\bar{x}^2}{\bar{x}}} \left( \frac{\bar{x}}{\bar{x}} \right)$$

$$\left[ a = \frac{\bar{x} \cdot \bar{y} - \bar{xy}}{(\bar{x})^2 - \bar{x}^2}, b = \bar{y} - a\bar{x} \right]$$

Para não limitarmos nossos coeficientes a e b, vamos chamá-los de  $\beta_1$  e  $\beta_0$ , respectivamente. Logo, a função da reta terá a forma

$$\hat{y} = \beta_1 \cdot x + \beta_0$$

ou seja:

$$\beta_0 = \bar{y} - a\bar{x} \quad \text{e} \quad \beta_1 = \frac{\bar{x} \cdot \bar{y} - \bar{xy}}{(\bar{x})^2 - \bar{x}^2}$$

Para efeitos de visualização ao multiplicarmos  $\beta_1$  por -1, temos:

$$\beta_1 = \frac{\bar{xy} - \bar{x} \cdot \bar{y}}{\bar{x}^2 - (\bar{x})^2}, \text{ que equivale a:}$$

$$\beta_1 = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

b) As suposições feitas sobre os erros em termos de distribuição, valor esperado e variância são: os resíduos seguem uma distribuição normal; os erros são independentes entre si, ou seja, a covariância entre os erros é 0; o modelo é linear nos parâmetros; e a homocedasticidade, que consiste em considerar que os erros além de seguirem uma distribuição normal, seguem a mesma normal ao longo de toda reta.; As suposições são importantes para que seja possível efetuar o teste de Hipóteses com os erros da equação, para que o modelo torne-se ainda mais confiável.

c) O teste de hipótese efetuado na regressão simples consiste em:

Hipótese nula:  $\beta_1 = 0$

Hipótese alternativa:  $\beta_1 \neq 0$

Este teste consiste em: ao rejeitarmos a hipótese nula, afirmamos que x e y possuem uma relação, ou seja, faz sentido fazermos uma regressão linear com as variáveis. Ao não rejeitarmos a hipótese nula, o contrário ocorre, comprovamos que não existe relação entre x e y.

d) Sim. Quanto ao modelo conforme descrito na equação, teríamos a aparição de mais um  $\beta_2.X_2$ ; Quanto as suposições do modelo, teríamos que acrescentar um eixo no gráfico da regressão, pois teríamos y dependendo de 2 variáveis ( $X_1$  e  $X_2$ ); Quanto ao teste de hipótese, teríamos que efetuar o teste de hipótese duas vezes, sendo uma com o beta de y e  $X_1$  e outra com o beta de y e  $X_2$ .