

## Ciência dos Dados Projeto 3 Parte 2

### Desenvolvimento Teórico da Técnica de Regressão

A seguir segue uma demonstração da dedução do método dos mínimos quadrados.

#### Método dos Mínimos Quadrados

• Representação de uma observação

$$L = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$[y_i - (\beta_0 + \beta_1 x_i)]^2 = (y_i - \beta_0 - \beta_1 x_i)^2$$

//

$$\left. \frac{d\epsilon}{d\hat{\beta}_0, \hat{\beta}_1} \right|_{\hat{\beta}_0, \hat{\beta}_1} = 2(y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^{2-1}(-1) + 2(y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^{2-1}(-1) + \dots + 2(y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^{2-1}(-1)$$

$$\left. \frac{d\epsilon}{d\hat{\beta}_0} \right|_{\hat{\beta}_0, \hat{\beta}_1} = (-2) \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \implies \boxed{n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i} \quad +$$

$$\left. \frac{d\epsilon}{d\hat{\beta}_1} \right|_{\hat{\beta}_0, \hat{\beta}_1} = (-2) \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 \implies \boxed{\hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n y_i x_i} \quad +$$

• Dividindo 1 por n:  $\hat{\beta}_0 + \hat{\beta}_1 \bar{x} = \bar{y} \implies \boxed{\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}} \quad +$

• Substituindo 3 em 2 e dividindo por  $n$ :

$$\begin{aligned} * \Rightarrow \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i - \underbrace{\hat{\beta}_0}_{3} \sum_{i=1}^n x_i \\ &= \sum_{i=1}^n x_i y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) \sum_{i=1}^n x_i \\ &= \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} + n \hat{\beta}_1 \bar{x}^2 \end{aligned}$$

$$\hat{\beta}_1 \sum_{i=1}^n x_i^2 - n \hat{\beta}_1 \bar{x}^2 = \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}$$

$$\hat{\beta}_1 \left( \sum_{i=1}^n x_i^2 - n \bar{x}^2 \right) = \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Sobre a distribuição dos erros podemos supor que este respeita uma distribuição normal de média zero e variância constante ou seja os dados regredidos encontram-se dispersos de forma homogênea em torno da reta de regressão, apresentando homocedasticidade. Na prática essas suposições podem ser checadas a partir da análise das funções de distribuição acumulada (fda) e da própria distribuição normal dos erros estabelecendo hipóteses e valores limites para que esta seja aceita.

Para uma regressão simples o teste de hipótese seria realizado para a hipótese nula  $\beta_1 = 0$ , dessa forma se a hipótese nula após o teste fosse

rejeitada poderíamos concluir que não existe relação linear entre as duas variáveis analisadas.

Além de regressões simples que estabelecem a relação entre uma variável explicativa e uma variável resposta, as regressões podem demonstrar a relação entre mais de duas variáveis sendo assim classificada como regressão múltipla que resultaria em uma equação mais complexa que envolvesse as outras variáveis de modo que a nova equação da reta seria:  $y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n + \varepsilon_i$ . E para verificar a relação entre as variáveis deveriam ser feitos testes de hipóteses para cada uma das variáveis explicativas (n testes).