

Cálculo Matricial

① Gradiente da constante é zero $\rightarrow f(\vec{x}) = Cte \Rightarrow \nabla f = \vec{0}$

② $f(\vec{x}) = \underbrace{\vec{x}^T}_{1 \times n} \cdot \underbrace{\vec{a}}_{n \times 1} \Rightarrow \nabla f = \vec{a}$
vetor constante

③ $f(\vec{x}) = \underbrace{\vec{x}^T}_{1 \times n} \underbrace{A}_{n \times n} \underbrace{\vec{x}}_{n \times 1} \Rightarrow \nabla f = (A + A^T) \vec{x}$
se $A = A^T \Rightarrow \nabla f = 2 \cdot A \vec{x}$

Note que: $\hat{y}^T y = y^T \hat{y}$
pois: $\underbrace{\hat{y}^T}_{1 \times m} y$ é matriz 1×1
 $\underbrace{y}_{m \times 1}$

Logo $(\hat{y}^T y) = (\hat{y}^T y)^T$ pois matriz 1×1 transposta é ela mesma
 $\hat{y}^T y = y^T \hat{y}$

Coisas de matriz: $(A+B)^T = A^T + B^T$
 $(AB)^T = B^T A^T$

Aplicando as regras do cálculo matricial ao MSE, temos:

$$\nabla \text{MSE} = \frac{1}{m} (2X^T X \theta - 2X^T y)$$

$$\nabla \text{MSE} = \vec{0} \Rightarrow \frac{1}{m} (2X^T X \theta - 2X^T y) = \vec{0} \Rightarrow$$

$$\Rightarrow 2X^T X \theta = 2X^T y \Rightarrow \boxed{\theta_{\text{opt}} = (X^T X)^{-1} X^T y}$$

Equações normal

Quanto custa calcular θ_{opt} ?

$$\theta_{\text{opt}} = \underbrace{(X^T X)^{-1}}_{O(mn^2)} \cdot \underbrace{X^T y}_{O(mn)} \xrightarrow{O(n^3)}$$

Modelo linear para regressão

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

notação matricial:

$$\hat{y} = X \theta, \text{ onde: } \hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1n} \\ 1 & X_{21} & X_{22} & \dots & X_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{m1} & X_{m2} & \dots & X_{mn} \end{bmatrix}$$

$$MSE = \frac{1}{m} (X\theta - y)^T (X\theta - y) = \frac{1}{m} (\theta^T X^T X \theta - 2\theta^T X^T y + y^T y)$$

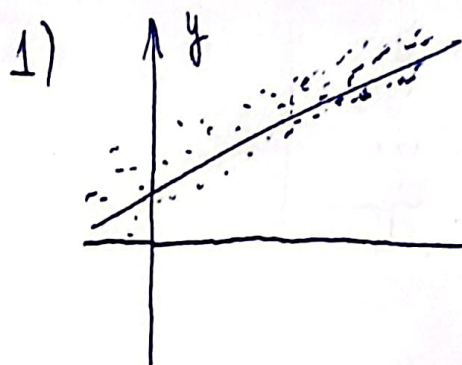
$$\nabla MSE = \frac{1}{m} (2X^T X \theta - 2X^T y)$$

Equação normal:

$$\nabla MSE = 0 \Rightarrow \frac{1}{m} (2X^T X \theta - 2X^T y) = \vec{0}$$

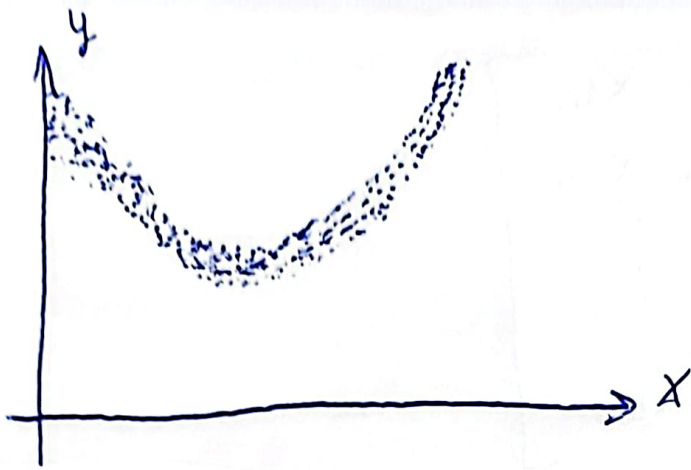
$$\theta = (X^T X)^{-1} \cdot X^T y$$

Regressão linear + feature engineering:



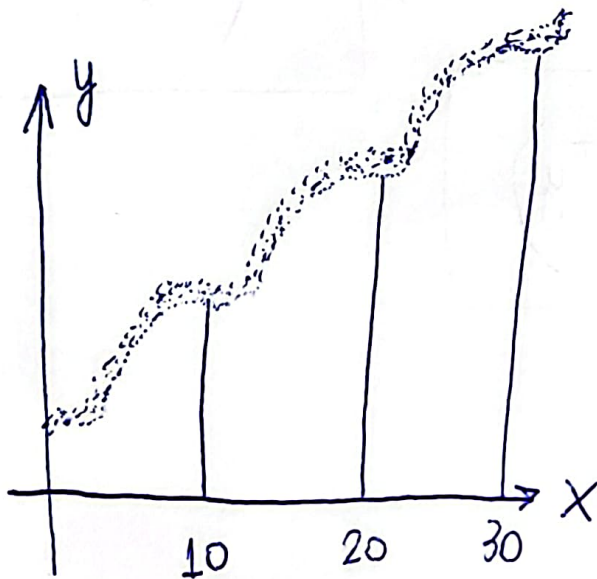
$$\hat{y} = \theta_0 + \theta_1 x$$

2)



$$\hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2$$

3)



$$\hat{y} = \theta_0 + \theta_1 x + \theta_2 \cdot \cos\left[\left(\frac{\pi}{5}\right)(x)\right] + \theta_3 \cdot \sin$$

$$\left[\left(\frac{\pi}{5}\right)(x)\right]$$

notação matricial:

$$1) \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}; X = \begin{bmatrix} 1 & x_{11} \\ 1 & x_{21} \\ \vdots & \vdots \\ 1 & x_{m1} \end{bmatrix}; y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$2) \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}; X = \begin{bmatrix} 1 & X_{11} & X_{11}^2 \\ 1 & X_{21} & X_{21}^2 \\ \dots & \dots & \dots \\ 1 & X_{m1} & X_{m1}^2 \end{bmatrix}$$

$$3) \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}; X = \begin{bmatrix} 1 & X_{11} & \cos(\dots) & \sin(\dots) \\ 1 & X_{21} & \cos(\dots) & \sin(\dots) \\ \dots & \dots & \dots & \dots \\ 1 & X_{m1} & \cos(\dots) & \sin(\dots) \end{bmatrix}$$

Intuição geométrica:

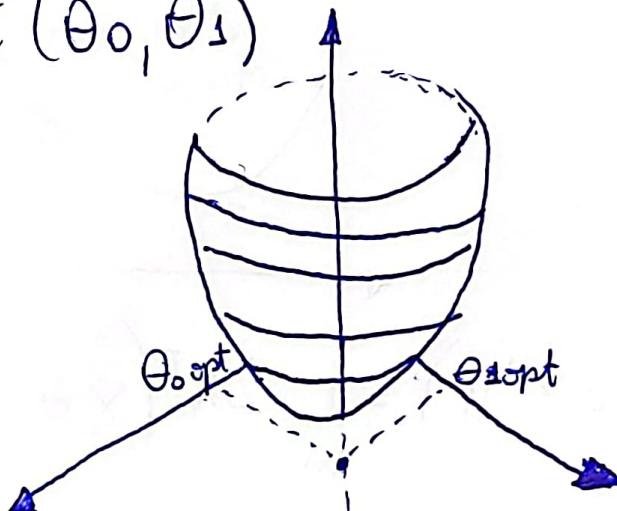
$$MSE(\theta) = \frac{1}{m} (\theta^T X^T X \theta - 2 \theta^T X^T y + y^T y)$$

Ata de ~~de~~ $MSE(\theta)$ é um parabolóide em θ

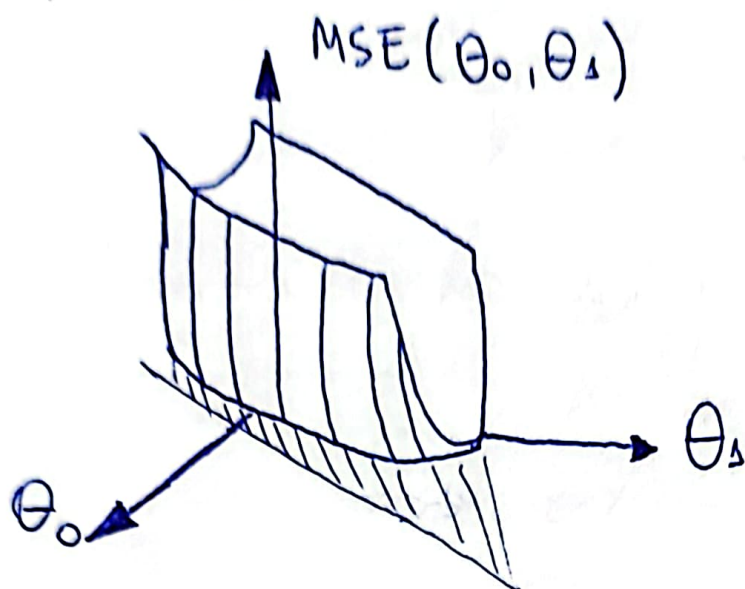
Exemplo:

$$\hat{y} = \theta_0 + \theta_1 X$$

$$MSE(\theta_0, \theta_1)$$



Exemplo degenerado:



Gradient Descent

1) Chuta $\theta^{(0)}$ ← iteração
 $i=0$

2) Enquanto não "acabou": $\theta^{(i+1)} = \theta^{(i)} - \eta \nabla \text{MSE}(\theta^{(i)})$
↳ taxa de aprendizado

3) Retorna o último θ

Lembrando que:

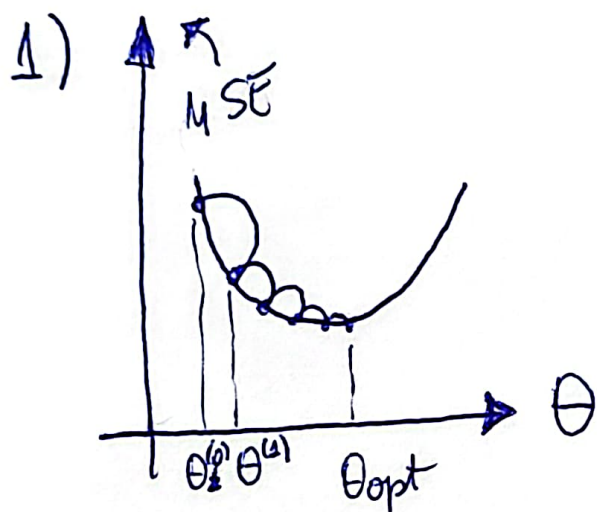
$$\nabla \text{MSE} = \frac{1}{m} (2X^T X \theta - 2X^T y) = \frac{1}{m} 2X^T (X\theta - y)$$

Critério de parada

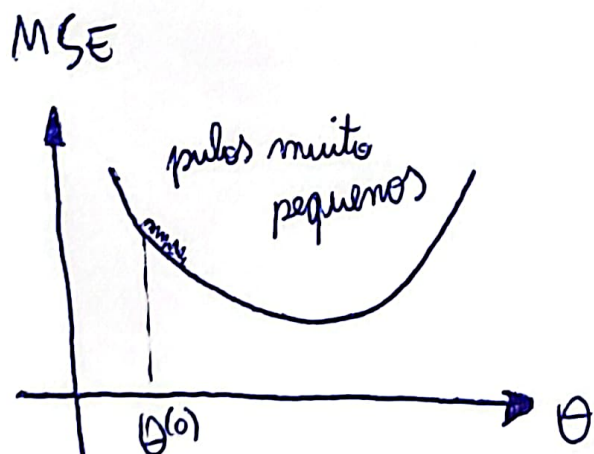
Se $\|\theta^{(i+1)} - \theta^{(i)}\| < \text{tolerância}$ ou se $i > \text{máximo de iterações}$

pode parar!

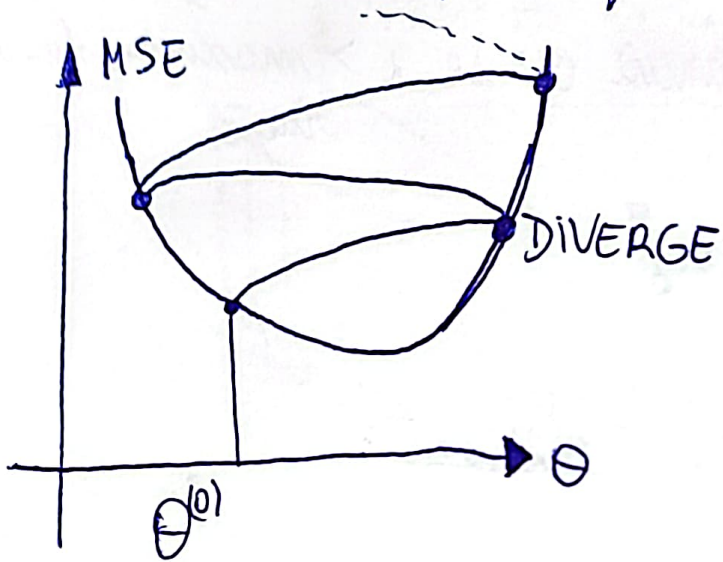
Intuição geométrica do gradient descent:



2) Taxa de aprendizagem muito pequena



3) Taxa de aprendizado muito grande



Modelos lineares

• conjunto de dados:

$$D = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)\}$$

onde: $\vec{x}_i = (f_1, f_2, \dots, f_n), f_k \in \mathbb{R}$

$y_i \in \mathbb{R}$ (regressão)

ou

$y_i \in \{\text{conjunto de classes}\}$ (classificação)

• notação matricial:

$$X_i = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{bmatrix} \quad (\text{matriz-coluna})$$

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_m^T \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$X \in M_{m,n}(\mathbb{R})$$

modelo: $h_{\theta}: \mathbb{R}^n \rightarrow \mathbb{R}$ (regressão)

↳ parâmetros treináveis { classes } (classificação)

denota previsão do modelo

↑

1

$$\hat{y} = h_{\theta}(x)$$

$$\rightarrow \hat{y} =$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix}$$

→ matriz-coluna (vetor) das previsões

erro de regressão

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

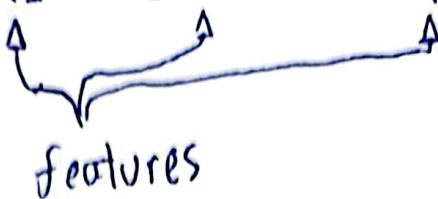
$$\epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_m \end{bmatrix} = \begin{bmatrix} \hat{y}_1 - y_1 \\ \hat{y}_2 - y_2 \\ \vdots \\ \hat{y}_m - y_m \end{bmatrix} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$\boxed{\epsilon = \hat{y} - y}$$

erro médio quadrático (MSE)

$$MSE = \frac{1}{m} \sum_{i=1}^m \epsilon_i^2 = \boxed{\frac{1}{m} \epsilon^T \cdot \epsilon}$$

Agora, o modelo linear

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + x_n \cdot \theta_n$$


features

- Treinar o modelo:

• Descobrir $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$ que minimiza o MSE

Um modelo de ML:

- 1) Dados
- 2) Função de erro
- 3) Modelo
- 4) Algoritmo de treinamento

Modelo Linear na forma matricial

Defina: $x_i' = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

Com isso: $\hat{y} = \theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \dots + \theta_n x_n$

$$\boxed{\hat{y} = x'^T \cdot \theta}$$

Portanto, defina:

$$X = \begin{bmatrix} 1 & \text{---} & X_1^T & \text{---} \\ 1 & \text{---} & X_2^T & \text{---} \\ \vdots & & \vdots & \\ 1 & \text{---} & X_m^T & \text{---} \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} = \begin{bmatrix} 1 & \text{---} & X_1^T & \text{---} \\ 1 & \text{---} & X_2^T & \text{---} \\ \dots & & \dots & \\ 1 & \text{---} & X_m^T & \text{---} \end{bmatrix} \cdot \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} = X \cdot \theta$$

Juntando as partes:

$$\varepsilon = \hat{y} - y; \quad \hat{y} = X\theta; \quad MSE = \frac{1}{m} \cdot \varepsilon^T \cdot \varepsilon$$

Logo:

$$MSE = \frac{1}{m} \cdot (\hat{y}^T \hat{y} - 2\hat{y}^T y + y^T y) = \text{~~scribbles~~}$$

$$= \frac{1}{m} ((X\theta)^T (X\theta) - 2(X\theta)^T y + y^T y) = \boxed{\frac{1}{m} (\theta^T X^T X \theta - 2\theta^T X^T y + y^T y)}$$

$$\nabla MSE = \begin{bmatrix} \frac{\partial MSE}{\partial \theta_0} \\ \frac{\partial MSE}{\partial \theta_1} \\ \frac{\partial MSE}{\partial \theta_2} \\ \vdots \\ \frac{\partial MSE}{\partial \theta_n} \end{bmatrix} = 0$$