

Modelo linear:

①

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

Nosso objetivo: achar os θ 's tal que nosso erro de treinamento seja o menor possível.

Conjunto de treinamento: X, y

$$X = \begin{bmatrix} | & | & \dots & | \\ x_1 & x_2 & \dots & x_n \\ | & | & \dots & | \end{bmatrix} \quad y = \begin{bmatrix} | \\ y \\ | \end{bmatrix}$$

~~exemplo 1, feature 1~~

$$X = \begin{bmatrix} \textcircled{x_{11}} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & \textcircled{x_{mn}} \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

→ exemplo 1, feature 1

→ exemplo m, feature n

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

previsão
p/
i-ésima
amostra
de
treinamento

(2)

IDEIA: Quero escrever as previsões p/ todas as amostras de treinamento de uma forma bem compacta, com matrizes.

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \vdots \\ \hat{y}^{(m)} \end{bmatrix} = \begin{bmatrix} \theta_0 + \theta_1 x_1^{(1)} + \dots + \theta_n x_n^{(1)} \\ \theta_0 + \theta_1 x_1^{(2)} + \dots + \theta_n x_n^{(2)} \\ \vdots \\ \theta_0 + \theta_1 x_1^{(m)} + \dots + \theta_n x_n^{(m)} \end{bmatrix}$$

$$= \begin{bmatrix} \theta_0 \cdot 1 + \theta_1 x_1^{(1)} + \dots + \theta_n x_n^{(1)} \\ \theta_0 \cdot 1 + \theta_1 x_1^{(2)} + \dots + \theta_n x_n^{(2)} \\ \vdots \\ \theta_0 \cdot 1 + \theta_1 x_1^{(m)} + \dots + \theta_n x_n^{(m)} \end{bmatrix}$$

$$= \underbrace{\begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix}}_{\substack{X \\ \sim \\ \uparrow \\ \text{"X's" aumentada}}} \underbrace{\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}}_{\theta}$$

⇒ todas as
previsões :

$$\hat{\mathbf{y}} = \mathbf{X} \boldsymbol{\theta}$$

Erro de treinamento

(3)

$$\varepsilon^{(i)} = \hat{y}^{(i)} - y^{(i)}$$

$$\mathcal{E} = \begin{bmatrix} \varepsilon^{(1)} \\ \varepsilon^{(2)} \\ \vdots \\ \varepsilon^{(m)} \end{bmatrix} = \begin{bmatrix} \hat{y}^{(1)} - y^{(1)} \\ \hat{y}^{(2)} - y^{(2)} \\ \vdots \\ \hat{y}^{(m)} - y^{(m)} \end{bmatrix} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \vdots \\ \hat{y}^{(m)} \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\Rightarrow \mathcal{E} = \hat{y} - y$$

$$SSE = [\varepsilon^{(1)}]^2 + [\varepsilon^{(2)}]^2 + \dots + [\varepsilon^{(m)}]^2$$

$$= \begin{bmatrix} \varepsilon^{(1)} & \varepsilon^{(2)} & \dots & \varepsilon^{(m)} \end{bmatrix} \begin{bmatrix} \varepsilon^{(1)} \\ \varepsilon^{(2)} \\ \vdots \\ \varepsilon^{(m)} \end{bmatrix}$$

$$= \mathcal{E}^T \cdot \mathcal{E} = (\hat{y} - y)^T (\hat{y} - y)$$

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

$$= ([\tilde{X}\theta]^T - y^T)(\tilde{X}\theta - y)$$

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

$$SSE = \theta^T \tilde{X}^T \tilde{X} \theta - 2\theta^T \tilde{X}^T y + y^T y$$

$$\frac{\partial SSE}{\partial \theta} = \nabla_{\theta} SSE = 2\tilde{X}^T \tilde{X} \theta - 2\tilde{X}^T y = 0$$

$$\Rightarrow 2\tilde{X}^T \tilde{X} \theta = 2\tilde{X}^T y \Rightarrow (\tilde{X}^T \tilde{X}) \theta = \tilde{X}^T y$$

$$\Rightarrow \theta^* = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T y$$

Equações
normais!

$$y = ax^2 + bx + c$$

$$\frac{\partial y}{\partial x} = 2ax + b$$

(4)

Com regularização Ridge:

$$\text{Perda} = \underline{SSE} + \alpha \sum_{j=0}^n \theta_j^2$$

$$\Rightarrow \text{Perda} = 2\theta^T (\tilde{X}^T \tilde{X} + \alpha I) \theta - 2\tilde{X}^T y + y^T y$$

$$\frac{\partial \text{Perda}}{\partial \theta} = 0 \Rightarrow 2(\tilde{X}^T \tilde{X} + \alpha I) \theta - 2\tilde{X}^T y = 0$$

$$\Rightarrow \theta^* = (\tilde{X}^T \tilde{X} + \alpha I)^{-1} \tilde{X}^T y$$

Equação normal da reg. Ridge.

$$\text{SVD: } X = U S V^T$$

U, V unitárias

$$S = \begin{bmatrix} \Sigma \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} \text{bloco} \\ \text{diagonal} \\ \text{não neg.} \\ \text{bloco} \\ \text{zeros} \end{matrix}$$

$$\text{Pois: } \tilde{X} = U S V^T$$

$$\tilde{X}^T \tilde{X} = V S^T \cancel{U^T U} S V$$

$$= V \begin{bmatrix} \Sigma^T & 0^T \end{bmatrix} \begin{bmatrix} \Sigma \\ \vdots \\ 0 \end{bmatrix} V = V \Sigma^2 V^T$$

~~Bezeledgers~~

$$(\tilde{X}^T \tilde{X})^{-1} = (V \Sigma^2 V^T)^{-1} = V \Sigma^{-2} V^T$$

$$\begin{bmatrix} 1/\sigma_1^2 & & 0 \\ & 1/\sigma_2^2 & \\ 0 & & \ddots \\ & & & 1/\sigma_n^2 \end{bmatrix}$$

$$(ABC)^{-1} = C^{-1} B^{-1} A^{-1}$$

Se X cols colinear $\Rightarrow \sigma_n = 0$ $\Rightarrow \Sigma^{-2}$ não existe!

Já com regularizações Ridge: $\alpha > 0$

⑤

$$\begin{aligned} & (\underline{X}^T \underline{X} + \alpha \mathbf{I}) \\ &= (\underline{V} \underline{S}^T \cancel{\underline{U}^T} \underline{U} \underline{S} \underline{V}^T + \alpha \mathbf{I}) \\ &= (\underline{V} \underline{\Sigma}^2 \underline{V}^T + \alpha \mathbf{I}) \\ &= (\underline{V} \underline{\Sigma}^2 \underline{V}^T + \underline{V} (\alpha \mathbf{I}) \underline{V}^T) \\ &= \underline{V} (\underline{\Sigma}^2 + \alpha \mathbf{I}) \underline{V}^T \end{aligned}$$

$$(\underline{X}^T \underline{X} + \alpha \mathbf{I})^{-1} = [\underline{V} (\underline{\Sigma}^2 + \alpha \mathbf{I}) \underline{V}^T]^{-1}$$

$$= \underline{V} \underbrace{(\underline{\Sigma}^2 + \alpha \mathbf{I})^{-1}} \underline{V}^T$$

$$\downarrow$$
$$\begin{bmatrix} (\sigma_1^2 + \alpha) & & & \\ & (\sigma_2^2 + \alpha) & & \\ & & \ddots & \\ 0 & & & (\sigma_n^2 + \alpha) \end{bmatrix}^{-1}$$

$$= \begin{bmatrix} \frac{1}{(\sigma_1^2 + \alpha)} & & & \\ & \frac{1}{(\sigma_2^2 + \alpha)} & & \\ & & \ddots & \\ & & & \frac{1}{(\sigma_n^2 + \alpha)} \end{bmatrix}$$