

Regularização

Regularização é inimiga do overfitting

Mas o que é a regularização?

Regularização



Aplicar
"vencas"
pseudas

Um exemplo de "vencas pseudas":

- uma moeda, vários lançamentos
- qual valor de $p = \text{Prob}(\text{dar "cara"})$

lançamento	p , só dados	p , com vença
—	indefinido	$1/2$
cara	1	$\approx 1/2$
cara	1	$\approx 1/2$
cara	$2/3$	$\approx 1/2$
cara	$3/4$	$\approx 1/2$
cara	$3/5$	$\approx 1/2$

matematicamente:

H : n caras (heads)

T : n coroa (tails)

1. só dados:

$$p = \frac{H}{H+T}$$

2. com vença a priori:

$$p = \frac{(H + \alpha)}{(H + \alpha) + (T + \alpha)}$$

$\alpha > 0$

E como se antes de lançar a moeda, já tivéssemos $H = \alpha$, $T = \alpha$

Regularização em modelos lineares

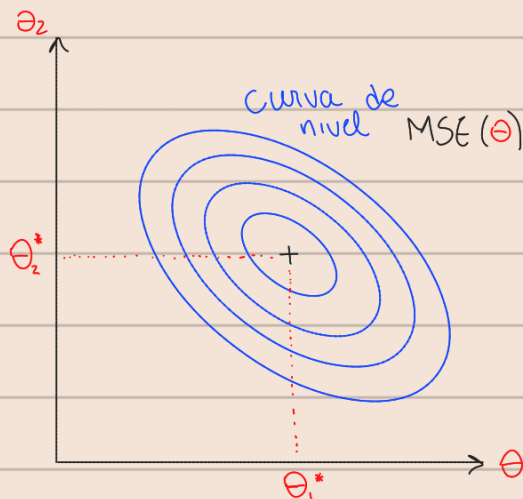
- sem regularização
- Ridge
- LASSO
- Elastic Net

* vol dados pequenos regularização ajuda

1) Sem regularização

$$L(\vec{\theta}) = \text{MSE}(\vec{\theta})$$

$$\underset{\vec{\theta}}{\text{argmin}} L(\vec{\theta}) \rightarrow \vec{\theta}_{\text{sem regularização}}^{\text{ótimo}}$$



+ ponto ótimo MSE

2) Regularização Ridge

* combate a colinearidade \rightarrow por causa da penalidade sobre o $\vec{\theta}^2$

* ã busca pelo ótimo, mas que eles sejam pequenos

\hookrightarrow vc quebra a possibilidade da minimização ser igual em vários casos

$$\vec{\theta}_{\text{Ridge}}^{\text{ótimo}} = \underset{\vec{\theta}}{\text{argmin}} \text{MSE}(\vec{\theta})$$

sujeito a $\|\vec{\theta}\|^2 \leq \beta$

norma l2

$$= \sqrt{\sum_{i=1}^n \|\theta_i\|^2}$$



multiplicadores de Lagrange

$$\vec{\theta}_{\text{Ridge}}^{\text{ótimo}} = \underset{\vec{\theta}}{\text{argmin}} \text{MSE}(\vec{\theta}) + \alpha \|\vec{\theta}\|^2$$

sem θ_0

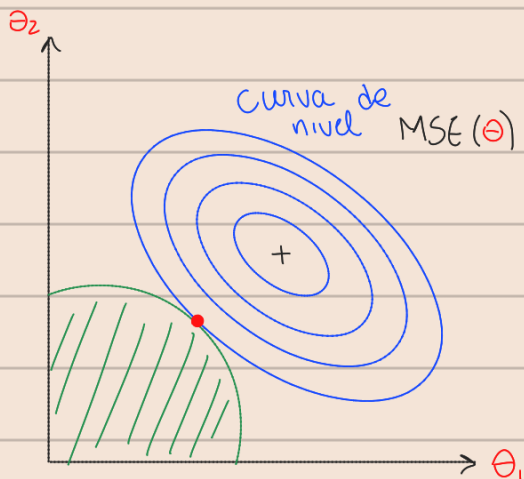
hiperparâmetro

Soluções:

- gradient descent
- equação normal

$$\vec{\theta}_{\text{Ridge}}^{\text{ótimo}} = (\underline{X}^T \underline{X} + \alpha \underline{I})^{-1} \underline{X}^T \underline{y}$$

* antes de uma regularização Ridge war Standard Scaler \rightarrow renome x_1 (milhões) não x_1 (micro) por exemplo



+ ponto ótimo MSE

Vantagens do Ridge

· combate a **colinearidade**

ex: $X_1 = X_2$

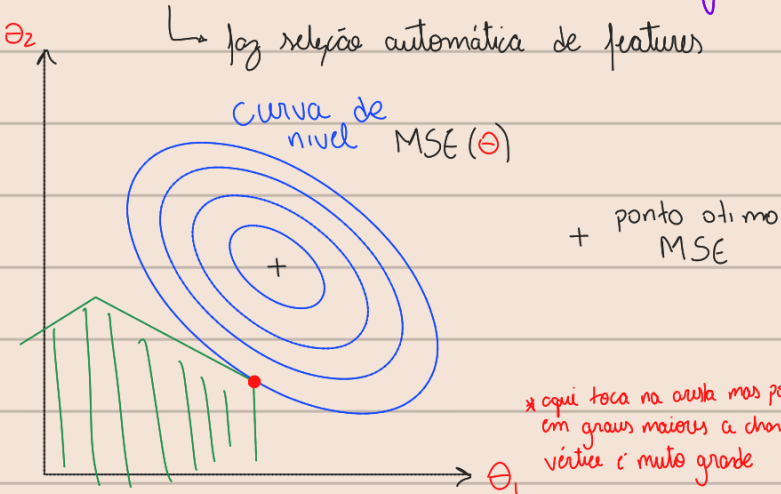
$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

1) Sem regularização, qualquer solução c/ $\theta_1 + \theta_2 = \theta_{12}^{\text{ótimo}}$ dá o mesmo modelo \rightarrow **ambiguo**

2) Ridge: $\theta_1 + \theta_2 = \frac{1}{2} \theta_{12}$ é a melhor opção

3) LASSO (least absolute shrinkage and selection operator)

\rightarrow faz seleção automática de features



$$\vec{\theta}_{\text{ótimo LASSO}} = \underset{\vec{\theta}}{\operatorname{argmin}} \operatorname{MSE}(\vec{\theta})$$

sujeito à $\|\vec{\theta}\|_1 \leq \beta$

$$\text{norma } l_1 = \sum_{i=1}^n |\theta_i|$$



$$\vec{\theta}_{\text{ótimo LASSO}} = \underset{\vec{\theta}}{\operatorname{argmin}} \operatorname{MSE}(\vec{\theta}) + \alpha \|\vec{\theta}\|_1$$

4) Reg. Elastic Net

$$\text{Loss} = \operatorname{MSE}(\vec{\theta}) + \alpha \cdot (1-\pi) \cdot \|\vec{\theta}\|_2^2 + \alpha \cdot \pi \cdot \|\vec{\theta}\|_1$$

α : penalidade

π : balançamento entre Ridge e LASSO

· combina vantagem dos dois modelos