

Um modelo de ML

- Dados
- Modelo
- Problema de otimização
- Algoritmo de otimização

Então o modelo linear é:

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

$$\hat{y} = 4 + 2x_1 + 3x_2$$
$$\hat{y} = 10 + 5x_1 + 15x_2$$

... etc

Então, o problema de otimização é minimizar o MSE em relação aos parâmetros $\vec{\theta} = (\theta_0, \theta_1, \theta_2)$

$$\vec{\theta}_{\text{opt}} = \underset{\vec{\theta}}{\text{argmin}} \text{MSE}(\vec{\theta}; X_{\text{train}}, y_{\text{train}})$$

O algoritmo de otimização é o procedimento usado para encontrar o ponto mínimo da função de perda

Exemplo: no modelo linear, temos alguns algoritmos de otimização:

1. Equação normal

1. Crie $\tilde{X} = \begin{bmatrix} \text{coluna de 1's} & X_{\text{train}} \end{bmatrix}$

↑
feature virtual

Um **modelo** é um conjunto de funções preditoras que obedecem à mesma regra de construção.

Exemplo: modelo linear para o problema de prever o **nível de açúcar** de uma banana baseado no seu **peso** e **comprimento**

$$y: \text{nvl açúcar} \quad x_1: \text{peso} \quad x_2: \text{comprimento}$$

O **problema de otimização** define um critério de escolha da **função preditiva** que melhor se ajusta a os **dados**

Exemplo: para o modelo linear é usual definir a **função de perda** como sendo o MSE

$$\text{MSE}(\vec{\theta}; X_{\text{train}}, y_{\text{train}}) = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

↳ $\theta_0 + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$

2. Resolva a equação normal

$$X^T X \vec{\theta} = X^T y_{\text{train}}$$

Se $m = n^\circ$ de exemplos

$n = n^\circ$ de features original (ou a virtual)

$$\Rightarrow X \in \mathbb{R}^{m \times (n+1)} \quad \vec{\theta} \in \mathbb{R}^{(n+1) \times 1}$$

$$y_{\text{train}} \in \mathbb{R}^{m \times 1}$$

$$\Rightarrow X X^T \in \mathbb{R}^{(n+1) \times (n+1)}$$

$$X^T y_{\text{train}} \in \mathbb{R}^{(n+1) \times 1}$$

$$X^T X \vec{\theta} \in \mathbb{R}^{(n+1) \times 1}$$

Custo de solução da equação normal?

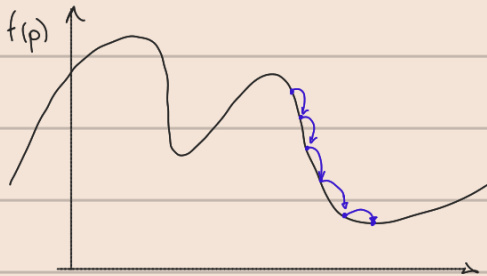
- custo construir $X^T X$ $O(n^2 m)$
- custo construir $X^T y_{\text{train}}$ $O(n m)$
- custo resolver sist linear $O(n^3)$

2. "Gradient descent"

Para X_{train} e y_{train} fixos, queremos minimizar

$$L(\vec{\theta}) = \text{MSE}(\vec{\theta}; X_{\text{train}}; y_{\text{train}})$$

gradient descent



1. Chuta valores para $\vec{\theta}$

$\vec{\theta} =$ chute inicial

2. Enquanto "não convergiu"

$$\vec{\theta} \leftarrow \vec{\theta} - \eta \nabla L(\vec{\theta})$$

novos $\vec{\theta}$

antigo $\vec{\theta}$

eta: tx de aprendizagem

$$\text{Precisamos de } \nabla L(\vec{\theta}) = \nabla \text{MSE}(\vec{\theta}; X_{\text{train}}, y_{\text{train}})$$

Felizmente esse gradiente é fácil de calcular no modelo linear:

$$\nabla L(\vec{\theta}) = \frac{2}{m} (X^T X \vec{\theta} - X^T y_{\text{train}})$$

$O(n m)$

Conclusão: se n° features muito alto \rightarrow gradient descent

se n° features baixo \rightarrow eq. normal

"O inimigo do overfitting é a regularização"