

ADAM

RESUMO TEÓRICO

Hiperparâmetros:

α : learning rate

β_1 : 0.9

β_2 : 0.999

ϵ : 10^{-8}

Variáveis:

t : época

g : vetor de gradientes

m : vetor de momentos

v : vetor de velocidades

w : vetor de pesos

Algoritmo (para cada época):

$$t := t + 1$$

$$g_t := \frac{\partial E_{t-1}}{\partial w_{t-1}}$$

$$m_t := \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$$

$$v_t := \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot (g_t \odot g_t)$$

$$\hat{m}_t := \frac{m_t}{1 - \beta_1}$$

$$\hat{v}_t := \frac{v_t}{1 - \beta_2}$$

$$w_t := w_{t-1} - \alpha \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

O ajuste dos pesos deve ser adaptado para o caso de mini-batches.