

1. We consider a feedforward network

$$a^{[0]} = x, \quad z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}, \quad a^{[l]} = \sigma^{[l]}(z^{[l]}), \quad l = 1, \dots, L$$

Assume all hidden layers use ReLU activation and the last layer is linear.

Since the output is linear, $f^{[L]} = 1$.

For $l = L-1, L-2, \dots, 1$:

$$f^{[l]} = (W^{[l+1]})^T f^{[l+1]} \odot \mathbb{1}_{z^{[l+1]} > 0}, \quad \text{where } \mathbb{1}_{z^{[l+1]} > 0} \text{ is an indicator function that equals 1 if } z^{[l+1]} > 0, \text{ and 0 otherwise.}$$

The final result is $\nabla_x a^{[L]}(x) = (W^{[L]})^T f^{[L]}$

2. (I) Practical issue

What happens if weights are initialized poorly?

Why do we sometimes add regularization in training, and how does it mathematically change the gradient?

(II). Optimization

Why is stochastic gradient descent (SGD) usually preferred over full gradient descent?

How does the learning rate impact convergence, and how should it be chosen in practice?

