

Part VII: Learning from Data

VII.1: Deep Neural Networks

VII.2: Convolutional Neural Nets

VII.3: Backpropagation

VII.4: Hyperparameters

VII.5: Machine Learning

VII.1 | The Construction of Deep Neural Networks

A neural network is a model $F: \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^k$ that makes a prediction $\hat{y} \in \mathbb{R}^k$ based on the parameters $x \in \mathbb{R}^n$ and the feature vector $v \in \mathbb{R}^m$:

$$\hat{y} = F(x, v).$$

Given the training data,

$$(v_i, y_i) \in \mathbb{R}^m \times \mathbb{R}^k, \quad i=1, \dots, N,$$

where v_i are features and y_i are labels, we optimize the parameters x so that the prediction of the model satisfies

$$y_i \approx F(x, v_i), \quad i=1, \dots, N.$$

In particular, neural networks are constructed from the following components:

- ① Affine functions: $v \mapsto Av + b$
- ② Activation functions: $w \mapsto \sigma(w)$
- ③ Composition: $F = F_L \circ \dots \circ F_2 \circ F_1$

Note that the composition of affine functions is affine:

$$F_1(v) = A_1 v + b_1, \quad F_2(v) = A_2 v + b_2$$

$$\begin{aligned} \Rightarrow (F_2 \circ F_1)(v) &= F_2(F_1(v)) \\ &= F_2(A_1 v + b_1) \\ &= A_2(A_1 v + b_1) + b_2 \\ &= A_2 A_1 v + (A_2 b_1 + b_2) \\ &= A v + b, \end{aligned}$$

where $A = A_2 A_1$ and $b = A_2 b_1 + b_2$.

The key to the modeling power of neural networks is the activation function.

The most popular activation function is the rectified linear unit:

$$\text{ReLU}(w) = [w]_+ = \begin{cases} 0 & \text{if } w \leq 0 \\ w & \text{if } w > 0 \end{cases}$$

Some other activation functions are:

(1) Linear:

$$\sigma(w) = a^T w + b,$$

(2) Hyperbolic tangent:

$$\sigma(w) = \tanh(w) \in (-1, 1),$$

(3) Logistic / Sigmoid:

$$\sigma(w) = \frac{1}{1 + e^{-w}} \in (0, 1).$$

Historically, \tanh and sigmoid were preferred for their smoothness.

ReLU makes $v \mapsto F(x, v)$ a continuous piecewise linear function.