# MULTILAYER NETWORKS

general idea



input layer          hidden layers          output layer

network's forward pass:

$$\mathbb{R}^{h^0} \ni x^0 \longmapsto \underset{\underset{\mathbb{R}^{h^1 \times h^0}}{\omega}}{W^1} \cdot x^1 + \underset{\underset{\mathbb{R}^{h^1}}{\omega}}{b^1} =: z^1$$

$$x^1 := \sigma^1(z^1) \longmapsto \underbrace{W^2 x^1 + b^1}_{} =: z^2$$

$$x^2 := \sigma^2(z^2) \longmapsto \underbrace{W^3 x^2 + b^3}_{} =: z^3$$

$$x^3 = \sigma^3(z^3)$$

- each layer takes as input $x^{k-1}$ and applies the maps $x^{k-1} \mapsto z^k := W^k x^{k-1} + b^k \longrightarrow \sigma^k(z^k)$

  where $W^k \in \mathbb{R}^{h^k \times h^{k-1}}$, $b^k \in \mathbb{R}^{h^k}$, $\sigma^k : \mathbb{R}^k \to \mathbb{R}^k$

  $(z_i)_{1 \le i \le k} \mapsto (\bar{\sigma}^k(z_i))_{1 \le i \le k}$

  for some $\bar{\sigma}^k : \mathbb{R} \to \mathbb{R}$

- the output of an N-layer neural network is

$$x^N = \sigma(z^N) = \sigma(W^N x^{N-1} + b^N) = \dots = x^N(x^0)$$

this is what we usually called the hypothesis function $h(x^0) \equiv x^N(x^0)$

- given the training data $(x^{(i)}, y^{(i)})_{1 \le i \le M}$ with $x^{(i)} \in \mathbb{R}^{h^0}$ being the features and $y^{(i)} \in \mathbb{R}^{h^N}$ the "labels":

  - for classification we use the convention

  $$\text{class } i \longleftrightarrow \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix} \in \mathbb{R}^{h^N} \quad \begin{matrix} 0 \\ i\text{-th position} \\ 1 \\ \vdots \\ N \end{matrix}$$

  - for regression $y^{(i)} \in \mathbb{R}^{h^N}$ are the $h^N$ continuous values of the output that is desired

After the forward pass we want to compare the output of the network to the labels of our training data and update the weights $W^e$ and biases $b^e$ accordingly:

$$L = \frac{1}{M} \sum_{i=1}^{M} l_i \qquad \text{as loss function}$$

for some function $l_i$ depending on the $i$-th pair $(x^{(i)}, y^{(i)})$, for example $l_i = \frac{1}{2}(y^{(i)} - x^N(x^{(i)}))^2$

Gradient descent: let $p^k$ be a place holder for some parameter in the $k$-th layer, i.e., $W_{ij}^k$, $b_i^k$ for $1 \leq i \leq n^k$, $1 \leq j \leq n^{k-1}$

$$\frac{\partial L}{\partial p^k} = \frac{1}{M} \sum_{i=1}^{M} \frac{\partial \ell_i}{\partial p^k} \qquad \text{recall} \quad \ell_i = \ell(y^{(i)}, x^N(x^{(i)}))$$

↑ depending on all the weights $W^\ell$ and biases $b^\ell$

$$\frac{\partial \ell}{\partial p^k} = \underbrace{\sum_{i=1}^{n^N} \frac{\partial}{\partial x_i^N} \ell(y, x^N)}_{=: \Delta(x^N)} \cdot \frac{\partial x_i^N}{\partial p^k}$$

$$\frac{\partial x^N}{\partial p^k} = \frac{\partial [\mathcal{G}(z^N)]}{\partial p^k} = \sum_{i=1}^{n^N} \mathcal{G}^{N'}(z_i^N) \frac{\partial z_i^N}{\partial p^k}$$

$$= \sum_{i=1}^{n^N} \mathcal{G}^{N'}(z_i^N) \frac{\partial \left[ \sum_{j=1}^{n^{N-1}} W_{ij}^N x_j^{N-1} + b_i^N \right]}{\partial p^k}$$

$$= \sum_{i=1}^{n^N} \mathcal{G}^{N'}(z_i^N) \left[ \sum_{j=1}^{n^{N-1}} W_{ij}^N \frac{\partial x_j^{N-1}}{\partial p^k} + \sum_{j=1}^{n^{N-1}} \frac{\partial W_{ij}^N}{\partial p^k} x_j^{N-1} + \frac{\partial b_i^N}{\partial p^k} \right]$$

We note that

$$\frac{\partial x^\ell}{\partial p^k} = 0 \quad \text{for} \quad \ell \leq k$$

$$\frac{\partial W^\ell}{\partial p^k} = 0 = \frac{\partial b^\ell}{\partial p^k} \quad \text{for} \quad \ell \neq k$$

Hence, the last expression simplifies depending on layer number $k$:

if $k \leq N-1$ this has to be expanded

$$\frac{\partial x^N}{\partial p^k} = \mathcal{G}^{N'}(z^N) \odot W^N \cdot \frac{\partial x^{N-1}}{\partial p^k}$$

$$+ \mathcal{G}^{N'}(z^N) \odot \left[ \underbrace{\frac{\partial W^N}{\partial p^k} x^{N-1} + \frac{\partial b^N}{\partial p^k}}_{\text{if } k \neq N \text{ this term is zero}} \right]$$

for $k < N$:

$$\frac{\partial x^N}{\partial p^k} = \mathcal{G}^{N'}(z^N) \odot W^N \cdot \frac{\partial x^{N-1}}{\partial p^k}$$

$$= \mathcal{G}^{N'}(z^N) \odot W^N \cdot \left[ \mathcal{G}^{N-1'}(z^{N-1}) \odot W^{N-1} \cdot \frac{\partial x^{N-1}}{\partial p^k} \right.$$

$$\left. + \mathcal{G}^{N-1'}(z^{N-1}) \odot \left( \frac{\partial W^{N-1}}{\partial p^k} x^{N-2} + \frac{\partial b^{N-1}}{\partial p^k} \right) \right]$$

and, hence, we get the following structure:

$$\frac{\partial \ell}{\partial p^k} = \Delta(x^N) \cdot \left[ \mathcal{G}^{N'}(z^N) \odot W^N \right] \cdot \frac{\partial x^{N-1}}{\partial p^k}$$

$$= \Delta(x^N) \cdot \left[ \mathcal{G}^{N'}(z^N) \odot W^N \right] \cdot \left[ \mathcal{G}^{N-1'}(z^{N-1}) \odot W^{N-1} \right]$$

$$\cdot \{ \ldots \} \cdot \left[ \mathcal{G}^{k+1'}(z^{k+1}) \odot W^{k+1} \right]$$

$$\cdot \left[ \mathcal{G}^{k'}(z^k) \odot \left( \frac{\partial W^k}{\partial p^k} x^{k-1} + \frac{\partial b^k}{\partial p^k} \right) \right]$$

## An efficient update rule: Backpropagation

Algorithm: 1) Input $x^0 \leftarrow x^{(i)}$

2) Forward pass: $x^0 \mapsto z^1 \mapsto x^1 \mapsto \ldots \mapsto z^N \mapsto x^N$

3) Compute $\Delta(x^N) = \frac{\partial}{\partial x_j} \ell(y^{(i)}, x_j) \Big|_{x_j^0 = x_{j(x^{(i)})}^N}$

4) Compute backward pass

$$\Delta^N := \Delta(x^N)$$
$$\Delta^{k-1} := \Delta^k \cdot \left[ \mathcal{Z}^{k\prime}(z^k) \odot W^k \right]$$

5) Compute

$$\frac{\partial \ell_i}{\partial p^k} = \Delta^k \frac{\partial x^k}{\partial p^k}$$

$$= \Delta^k \left[ \mathcal{Z}^{k\prime}(z^k) \odot \left( \frac{\partial W^k}{\partial p^k} \cdot x^{k-1} + \frac{\partial b^k}{\partial p^k} \right) \right]$$

6) Compute $\ell_i$ for $i \in I$ (mini-batch $I \subseteq \{1 \ldots n\}$)

and average

$$\frac{\partial L_I}{\partial p^k} = \frac{1}{|I|} \sum_{i \in I} \frac{\partial \ell_i}{\partial p^k}$$

7) Update weights accordingly

$$W_{ij}^k \mapsto W_{ij}^k - \eta \frac{\partial L_I}{\partial W_{ij}^k}$$

$$b_i^k \mapsto b_i^k - \eta \frac{\partial L_I}{\partial b_i^k}$$

8) repeat for all mini-batches and epochs.

---

[HW] Derive this backpropagation algorithm from the KKT Condition for the Lagrangian formulation of:

$$\min_{W^k, b^k} \mathcal{L} = \frac{1}{n} \sum_{i=1}^{M} \ell(y^{(i)}, x^N)$$

under constraints $x^k = \mathcal{Z}(W^k x^{k-1} + b^k)$

1) formulate Lagrangian

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^{M} \ell(y^{(i)}, x^N) + \sum_{k=1}^{M} \sum_{k=1}^{N} \beta_i^k \cdot \left[ x^k - \mathcal{Z}(W^k x^{k-1} + b^k) \right]$$

2) $\frac{\partial \mathcal{L}}{\partial \beta_i}$ gives forward pass

3) $\frac{\partial \mathcal{L}}{\partial x_i}$ gives backward pass

4) $\frac{\partial \mathcal{L}}{\partial W^k}, \frac{\partial \mathcal{L}}{\partial b^k}$ give update rule