

Memo: Back-propagation

Tomomichi Sugihara

March 27, 2021

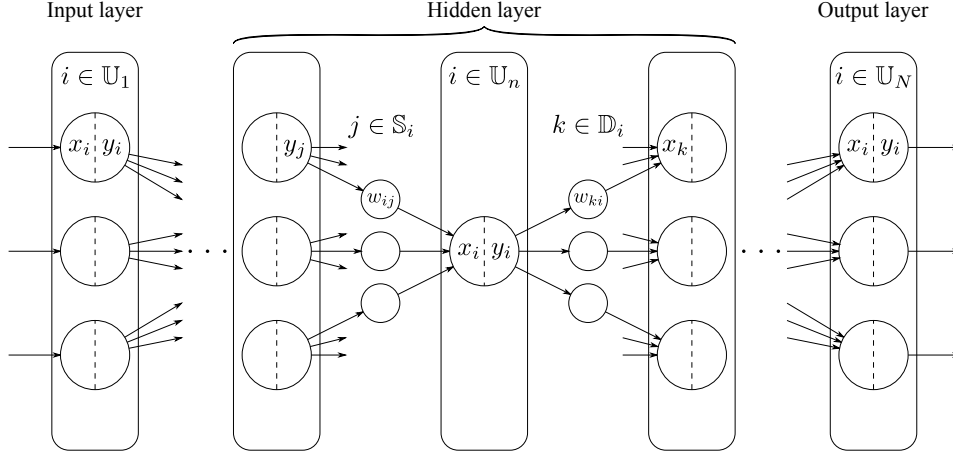


Figure 1: A non-recurrent multi-layered neural network

Let us consider a non-recurrent N -layered neural network illustrated in Fig. 1. The transmission of information in it is represented by the following equations:

$$x_i = \begin{cases} \text{given} & (\text{if } i \in \mathbb{U}_1) \\ \sum_{j \in \mathbb{S}_i} w_{ij} y_j + b_i & (\text{otherwise}) \end{cases} \quad (1)$$

$$y_i = \begin{cases} x_i & (\text{if } i \in \mathbb{U}_1) \\ f_i(x_i), & (\text{otherwise}) \end{cases} \quad (2)$$

where x_i is the input to the i th neural unit, \mathbb{U}_n for $n = 1, \dots, N$ is the set of indices of units in the n th layer, \mathbb{S}_i is the set of indices of upstream neural units connected to the i th unit, y_i is the output of the i th unit, w_{ij} is the weight on y_j that is input to the i th unit, b_i is the bias of the i th unit, and $f_i(\cdot)$ is the activation function associated with the i th unit. By propagating the output of each unit from the input layer to the output layer based on the above equations, the network transforms the input $\mathbf{x} = \{x_i | \forall i \in \mathbb{U}_1\}$ to the output $\mathbf{y} = \{y_i | \forall i \in \mathbb{U}_N\}$. Namely, it represents a multi-input-multi-output algebraic mapping $\mathcal{N}_{\{w_{ij}\}, \{b_i\}} : \mathbf{x} \mapsto \mathbf{y}$ or $\mathbf{y} = \mathbf{F}(\mathbf{x})$.

Provided an input \mathbf{x}^* paired with the desired output \mathbf{y}^* , the loss function of the network $E(\mathbf{x}^*, \mathbf{y}^*; \mathcal{N}_{\{w_{ij}\}, \{b_i\}})$ is defined. A typical choice of the function is the sum of squared errors as

$$E(\mathbf{x}^*, \mathbf{y}^*; \mathcal{N}_{\{w_{ij}\}, \{b_i\}}) \stackrel{\text{def}}{=} \frac{1}{2} \|\mathbf{F}(\mathbf{x}^*) - \mathbf{y}^*\|^2. \quad (3)$$

The back-propagation is one of the representative methods to train the network, *i.e.*, all the weights $\{w_{ij}\}$ and the biases $\{b_i\}$ so as to minimize the above loss function based on the following update rule:

$$w_{ij} \leftarrow w_{ij} - \eta \Delta w_{ij} \quad (4)$$

$$b_i \leftarrow b_i - \eta \Delta b_i, \quad (5)$$

where η is the learning rate. Δw_{ij} and Δb_i are decided to align the steepest descent direction as

$$\Delta w_{ij} = \frac{\partial E_{\mathbb{L}}}{\partial w_{ij}} \quad (6)$$

$$\Delta b_i = \frac{\partial E_{\mathbb{L}}}{\partial b_i}, \quad (7)$$

where $E_{\mathbb{L}}$ is the following loss function of a mini-batch

$$E_{\mathbb{L}} = \sum_{l \in \mathbb{L}} E(\mathbf{x}^{*(l)}, \mathbf{y}^{*(l)}; \mathcal{N}_{\{w_{ij}\}, \{b_i\}}), \quad (8)$$

\mathbb{L} is the set of identifiers of runs included in the mini-batch, $\mathbf{x}^{*(l)}$ is the input in the l th run to the network, and $\mathbf{y}^{*(l)}$ is the corresponding desired output of the network. Obviously,

$$\frac{\partial E_{\mathbb{L}}}{\partial w_{ij}} = \sum_{l \in \mathbb{L}} \frac{\partial E}{\partial w_{ij}}(\mathbf{x}^{*(l)}, \mathbf{y}^{*(l)}; \mathcal{N}_{\{w_{ij}\}, \{b_i\}}) \quad (9)$$

$$\frac{\partial E_{\mathbb{L}}}{\partial b_i} = \sum_{l \in \mathbb{L}} \frac{\partial E}{\partial b_i}(\mathbf{x}^{*(l)}, \mathbf{y}^{*(l)}; \mathcal{N}_{\{w_{ij}\}, \{b_i\}}), \quad (10)$$

and hence, the goal here is to find $\frac{\partial E}{\partial w_{ij}}$ and $\frac{\partial E}{\partial b_i}$.

w_{ij} and b_i only affect x_i directly, and thus,

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial x_i} \frac{\partial x_i}{\partial w_{ij}} = p_i y_j \quad (11)$$

$$\frac{\partial E}{\partial b_i} = \frac{\partial E}{\partial x_i} \frac{\partial x_i}{\partial b_i} = p_i, \quad (12)$$

where

$$p_i \stackrel{\text{def}}{=} \frac{\partial E}{\partial x_i}. \quad (13)$$

x_i affects y_i , and y_i affects $\{x_k | \forall k \in \mathbb{D}_i\}$, where \mathbb{D}_i is the set of indices of downstream units connected to the i th unit, if it belongs to a hidden layer. Hence,

$$p_i = \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial x_i} = \begin{cases} \frac{\partial E}{\partial y_i} f'_i(x_i) & (\text{if } i \in \mathbb{U}_N) \\ \left(\sum_{k \in \mathbb{D}_i} \frac{\partial E}{\partial x_k} \frac{\partial x_k}{\partial y_i} \right) \frac{\partial y_i}{\partial x_i} = \left(\sum_{k \in \mathbb{D}_i} p_k w_{ki} \right) f'_i(x_i). & (\text{otherwise}) \end{cases} \quad (14)$$

In the typical case of the sum of squared errors,

$$\frac{\partial E}{\partial y_i} = y_i - y_i^* \quad \text{for } \forall i \in \mathbb{U}_N. \quad (15)$$

$f'_i(x_i)$ depends on the definition of $f_i(x_i)$. If it is the sigmoid function,

$$f_i(x_i) = \frac{1}{1 + e^{-x_i}} \quad \Rightarrow \quad f'_i(x_i) = \frac{e^{-x_i}}{(1 + e^{-x_i})^2} = y_i(1 - y_i). \quad (16)$$

Or, if it is ReLU,

$$f_i(x_i) = \max\{x_i, 0\} \quad \Rightarrow \quad f'_i(x_i) = \begin{cases} 1 & (\text{if } x_i \geq 0) \\ 0 & (\text{otherwise}) \end{cases}. \quad (17)$$

Note that ReLU is not differentiable at $x_i = 0$ in the strict sense, although it is usually ignored.

A pseudocode of an algorithm for a neural network $\mathcal{N}_{\{w_{ij}\}, \{b_i\}}$ to train from data set $\left\{ (\mathbf{x}^{*(l)}, \mathbf{y}^{*(l)}) \mid \forall l \in \mathbb{L} \right\}$ and the loss function $E(\mathbf{x}^*, \mathbf{y}^*; \mathcal{N}_{\{w_{ij}\}, \{b_i\}})$ based on the above back-propagation is as follows.

Algorithm 1 TRAINNN($\mathcal{N}_{\{w_{ij}\}, \{b_i\}}, \{(\mathbf{x}^{*(l)}, \mathbf{y}^{*(l)}) \mid \forall l \in \mathbb{L}\}, \eta$)

Require: $\mathcal{N}_{\{w_{ij}\}, \{b_i\}} \cdots$ initialized

Ensure: $\mathcal{N}_{\{w_{ij}\}, \{b_i\}} \cdots$ updated

- 1: INITGRAD($\mathcal{N}_{\{w_{ij}\}, \{b_i\}}$) ▷ Initialize gradients of weights and biases
 - 2: **for** $l \in \mathbb{L}$ **do** ▷ Evaluate each run in a mini-batch
 - 3: BACKPROPAGATE($\mathcal{N}_{\{w_{ij}\}, \{b_i\}}, (\mathbf{x}^{*(l)}, \mathbf{y}^{*(l)})$)
 - 4: **end for**
 - 5: UPDATENN($\mathcal{N}_{\{w_{ij}\}, \{b_i\}}, \eta$) ▷ Update weights and biases of neural network
-

Algorithm 2 INITGRAD($\mathcal{N}_{\{w_{ij}\}, \{b_i\}}$)

- 1: **for** $i \in \mathbb{U}_2 \cup \cdots \cup \mathbb{U}_N$ **do**
 - 2: **for** $j \in \mathbb{S}_i$ **do**
 - 3: $\Delta w_{ij} \leftarrow 0$
 - 4: **end for**
 - 5: $\Delta b_i \leftarrow 0$
 - 6: **end for**
-

Algorithm 3 BACKPROPAGATE($\mathcal{N}_{\{w_{ij}\}, \{b_i\}}, (\mathbf{x}^{*(l)}, \mathbf{y}^{*(l)})$)

Ensure: $\{\Delta w_{ij}\}, \{\Delta b_i\}$

- 1: PROPAGATE($\mathcal{N}_{\{w_{ij}\}, \{b_i\}}, \mathbf{x}^{*(l)}$)
 - 2: INITPARAM($\mathcal{N}_{\{w_{ij}\}, \{b_i\}}$)
 - 3: **for** $i \in \mathbb{U}_N$ **do** ▷ for output layer
 - 4: $p_i \leftarrow \partial E / \partial y_i(\mathbf{x}^{*(l)}, \mathbf{y}^{*(l)}; \mathcal{N}_{\{w_{ij}\}, \{b_i\}})$ (dependent on the definition of E)
 - 5: **end for**
 - 6: **for** $n = N$ **downto** 2 **do** ▷ Backpropagation
 - 7: **for** $i \in \mathbb{U}_n$ **do**
 - 8: $p_i \leftarrow p_i v_i$
 - 9: **for** $j \in \mathbb{S}_i$ **do**
 - 10: $p_j \leftarrow p_j + p_i w_{ij}$
 - 11: $\Delta w_{ij} \leftarrow \Delta w_{ij} + p_i y_j$
 - 12: **end for**
 - 13: $\Delta d_i \leftarrow \Delta d_i + p_i$
 - 14: **end for**
 - 15: **end for**
-

Algorithm 4 INITPARAM($\mathcal{N}_{\{w_{ij}\}, \{b_i\}}$)

- 1: **for** $i \in \mathbb{U}_1 \cup \cdots \cup \mathbb{U}_N$ **do**
 - 2: $p_i \leftarrow 0$
 - 3: $v_i \leftarrow f'(x_i)$
 - 4: **end for**
-

Algorithm 5 UPDATENN($\mathcal{N}_{\{w_{ij}\}, \{b_i\}}, \eta$)

Require: $\{\Delta w_{ij}\}, \{\Delta b_i\} \cdots$ computed

Ensure: $\{w_{ij}\}, \{b_i\} \cdots$ updated

- 1: **for** $i \in \mathbb{U}_2 \cup \cdots \cup \mathbb{U}_N$ **do**
 - 2: $d_i \leftarrow d_i - \eta \Delta d_i$
 - 3: **for** $j \in \mathbb{S}_i$ **do**
 - 4: $w_{ij} \leftarrow w_{ij} - \eta \Delta w_{ij}$
 - 5: **end for**
 - 6: **end for**
-