

Neural Network: Theory and Practice

Jephian C.-H. Lin 林晉宏

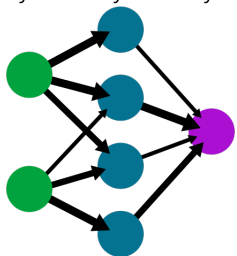
Department of Applied Mathematics, National Sun Yat-sen University

February 10, 2025

The 2nd Sizihwan Combinatorics Conference

What is an artificial neural network?

A simple neural network
input layer hidden layer output layer



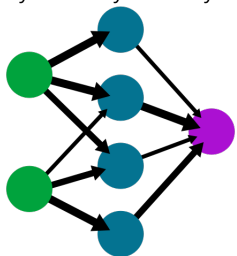
(Source: Wikipedia of Neural network)

keep improving itself

Not easy to find learning resources for mathematicians.

What is an artificial neural network?

A simple neural network
input layer hidden layer output layer

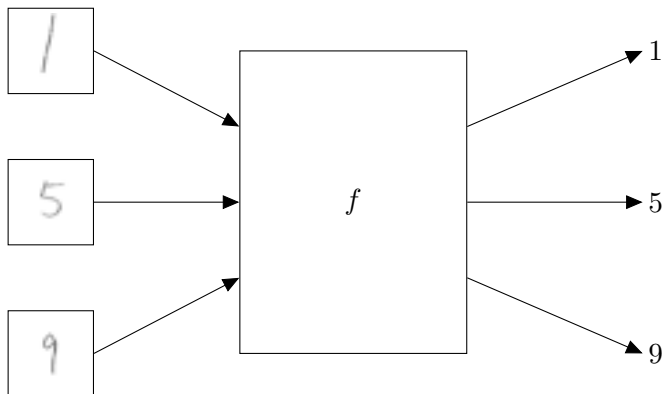


(Source: Wikipedia of Neural network)

keep improving itself

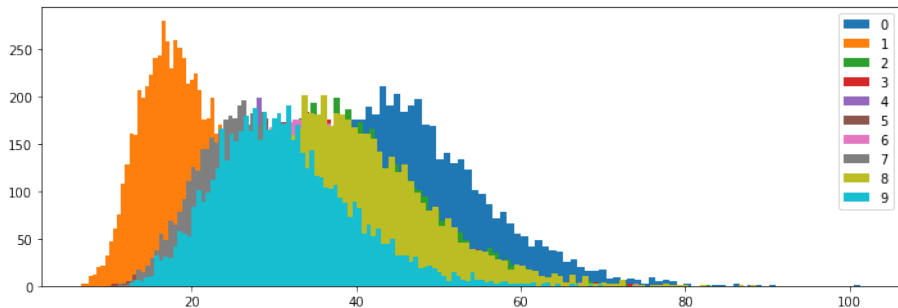
Not easy to find learning resources for mathematicians.

Classification task

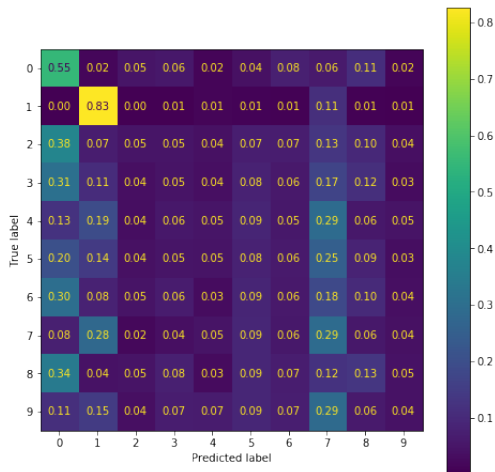


Function f

- $f(\mathbf{x}_i^\top) = 1$ (accuracy 10%)
- $f(\mathbf{x}_i^\top) = \mathbf{x}_i^\top \mathbf{1} \rightarrow$ guess the color by ink amount (accuracy 22%)



Confusion matrix of the ink-guessing strategy

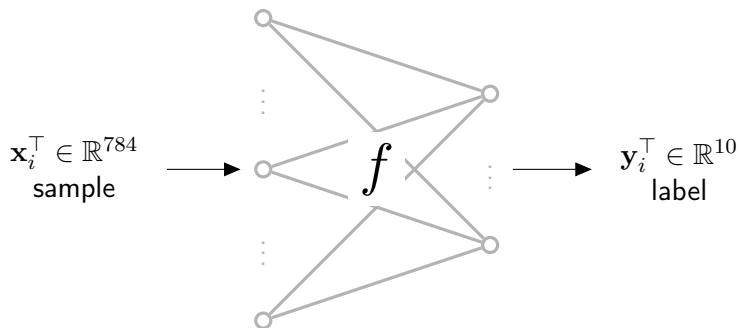


Output: One-hot encoding

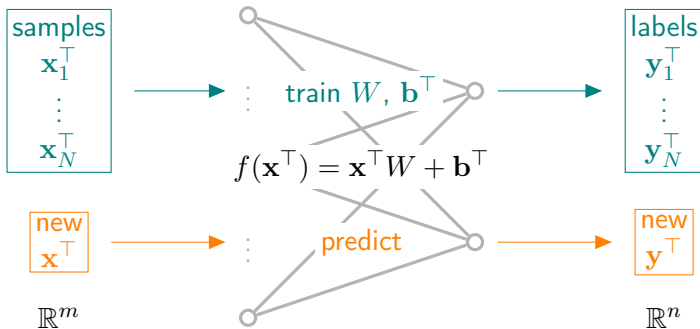
$$\begin{array}{r} 1 \rightarrow [0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0] \\ \vdots \qquad \qquad \qquad \qquad \qquad \qquad \vdots \\ 5 \rightarrow [0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0] \in \mathbb{R}^{10} \\ \vdots \qquad \qquad \qquad \qquad \qquad \qquad \vdots \\ 9 \rightarrow [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1] \end{array}$$

One-hot encoding ensures the labels of $0, \dots, 9$ are mutually equal-distanced.

A layer \sim a transformation



Linear model



Problem

Find $W \in \mathbb{R}^{m \times n}$ and $\mathbf{b}^\top \in \mathbb{R}^n$ that minimizes the *cost function*

$$\frac{1}{N} \sum_{i=1}^N \|f(\mathbf{x}_i^\top) - \mathbf{y}_i\|^2.$$

Solve the optimization problem — by projection

Let

$$X = \begin{bmatrix} - & \mathbf{x}_1 & - \\ & \vdots & \\ - & \mathbf{x}_N & - \end{bmatrix} \in \mathbb{R}^{N \times m} \text{ and } Y = \begin{bmatrix} - & \mathbf{y}_1 & - \\ & \vdots & \\ - & \mathbf{y}_N & - \end{bmatrix} \in \mathbb{R}^{N \times n}.$$

The cost function becomes

$$\frac{1}{N} \|XW + \mathbf{1}_N \mathbf{b}^\top - Y\|^2 = \frac{1}{N} \|\hat{X}\hat{W} - Y\|^2,$$

where

$$\hat{X} = [\mathbf{1}_N \quad X] \text{ and } \hat{W} = \begin{bmatrix} \mathbf{b}^\top \\ W \end{bmatrix}.$$

The minimum uniquely occurs at $\hat{W} = (\hat{X}^\top \hat{X})^{-1} \hat{X}^\top Y$.

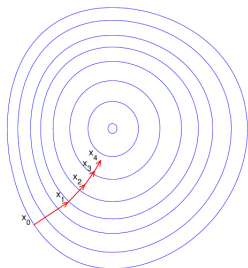
Memory consuming. Hard to update for new data.

Solve the optimization problem — by gradient descent

Input a scalar-valued function f and a **learning rate** $\alpha = 0.1$

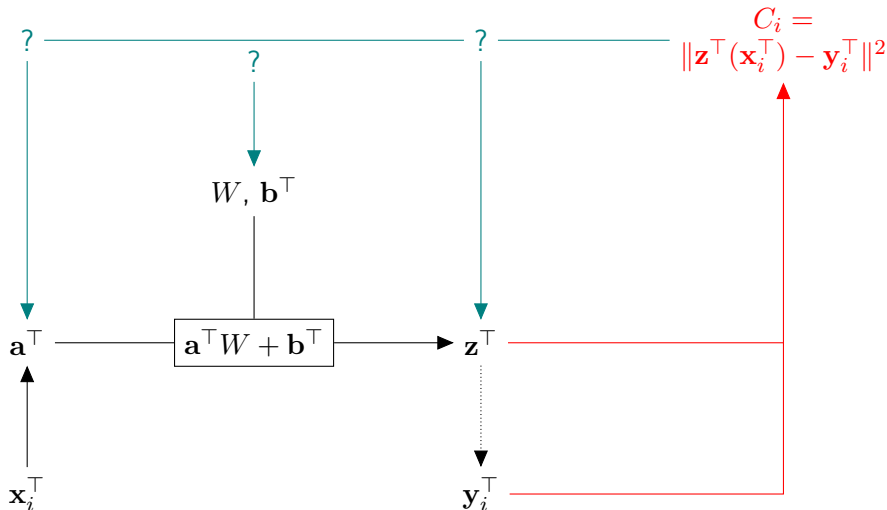
Output a point \mathbf{x} such that $f(\mathbf{x})$ is **locally** minimized, **hopefully**

- Steps**
- 1 Initiate \mathbf{x} randomly.
 - 2 Compute $\mathbf{v} = \nabla f(\mathbf{x})$.
 - 3 $\mathbf{x} \leftarrow \mathbf{x} - \alpha \mathbf{v}$.
 - 4 Repeat Step 2 \sim Step 3 several iterations.



(Source: Wikipedia of Gradient descent)

Linear model



Derivative and gradient

Derivative is a linear operator
sending a direction to the directional derivative.

$$\frac{\partial f}{\partial \mathbf{x}}(\Delta_{\mathbf{x}}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\Delta_{\mathbf{x}}) - f(\mathbf{x})}{t} \stackrel{\text{if } f \text{ scalar-valued}}{\text{gradient}} \nabla_{\mathbf{x}} f \cdot \Delta_{\mathbf{x}}.$$

Example

Let $C_i = \|\mathbf{z}^{\top}(\mathbf{x}_i^{\top}) - \mathbf{y}_i^{\top}\|$. Then

$$\begin{aligned} \frac{\partial C_i}{\partial \mathbf{z}^{\top}}(\Delta_{\mathbf{z}^{\top}}) &= \lim_{t \rightarrow 0} \frac{\|\mathbf{z}^{\top} + t\Delta_{\mathbf{z}^{\top}} - \mathbf{y}_i^{\top}\|^2 - \|\mathbf{z}^{\top} - \mathbf{y}_i^{\top}\|^2}{t} \\ &= \lim_{t \rightarrow 0} \frac{(\mathbf{z}^{\top} + t\Delta_{\mathbf{z}^{\top}} - \mathbf{y}_i^{\top}) \cdot (\mathbf{z}^{\top} + t\Delta_{\mathbf{z}^{\top}} - \mathbf{y}_i^{\top}) - \|\mathbf{z}^{\top} - \mathbf{y}_i^{\top}\|^2}{t} \\ &= 2(\mathbf{z}^{\top} - \mathbf{y}_i^{\top}) \cdot \Delta_{\mathbf{z}^{\top}}. \end{aligned}$$

Derivative and gradient

Derivative is a linear operator
sending a direction to the directional derivative.

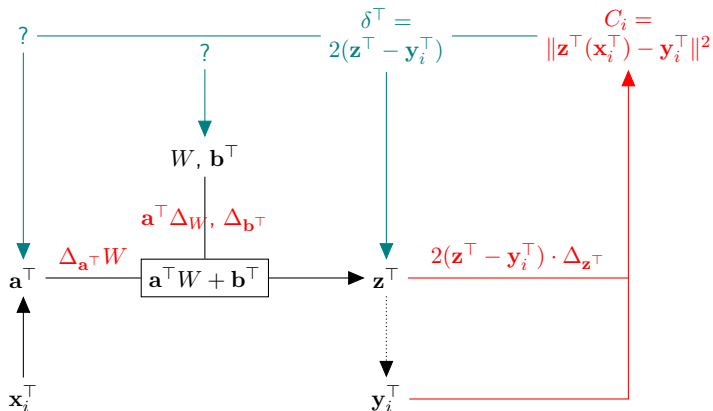
$$\frac{\partial f}{\partial \mathbf{x}}(\Delta_{\mathbf{x}}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\Delta_{\mathbf{x}}) - f(\mathbf{x})}{t} \stackrel{\text{if } f \text{ scalar-valued}}{=} \underbrace{\nabla_{\mathbf{x}} f}_{\text{gradient}} \cdot \Delta_{\mathbf{x}}.$$

Example

Let $f(\mathbf{a}^\top) = \mathbf{a}^\top W + \mathbf{b}^\top$. Then

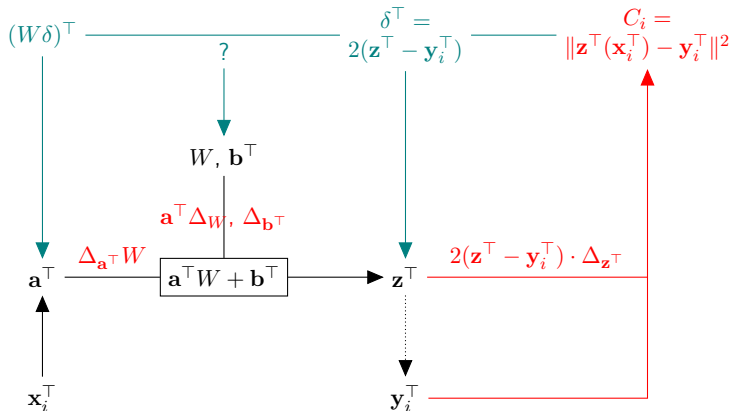
$$\frac{\partial f}{\partial \mathbf{a}^\top}(\Delta_{\mathbf{a}^\top}) = \Delta_{\mathbf{a}^\top} W, \quad \frac{\partial f}{\partial W}(\Delta_W) = \mathbf{a}^\top \Delta_W, \quad \frac{\partial f}{\partial \mathbf{b}^\top}(\Delta_{\mathbf{b}^\top}) = \Delta_{\mathbf{b}^\top}.$$

Linear model



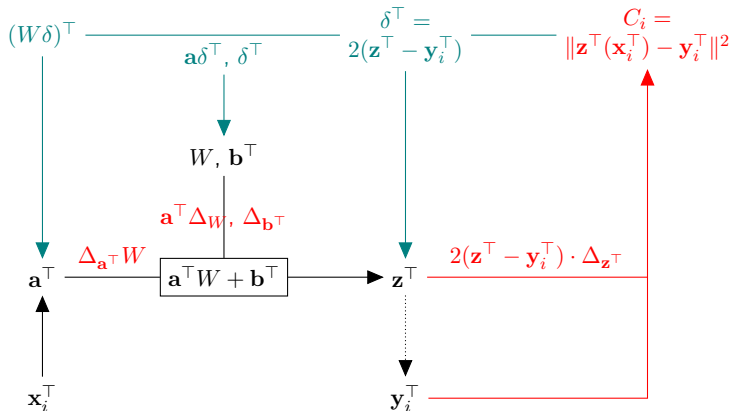
Linear model

$$\frac{\partial C_i}{\partial \mathbf{a}} = \Delta_{\mathbf{a}^\top} W \delta = (W \delta)^\top \cdot \Delta_{\mathbf{a}^\top}$$

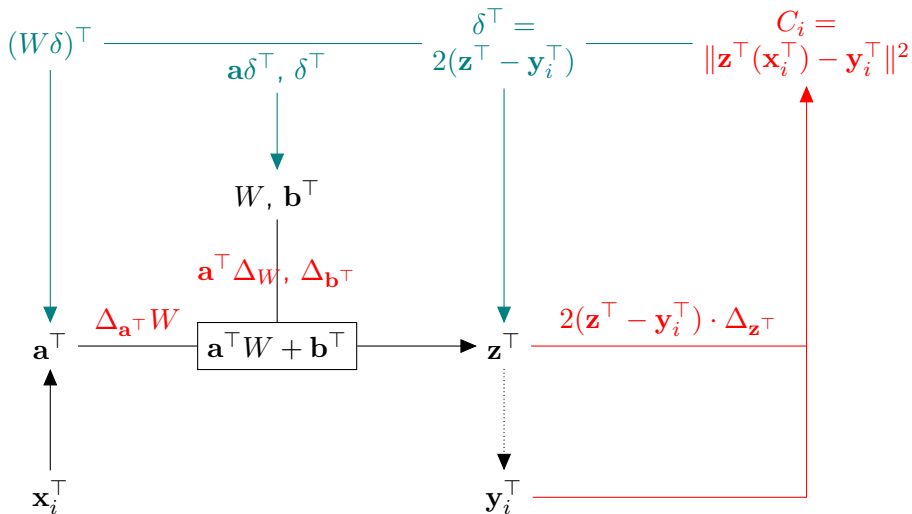


Linear model

$$\frac{\partial C_i}{\partial W} = \mathbf{a}^\top \Delta_W \delta = (\mathbf{a} \delta^\top) \cdot \Delta_W, \quad \frac{\partial C_i}{\partial \mathbf{b}^\top} = \delta^\top \cdot \Delta_{\mathbf{b}^\top}$$



Linear model

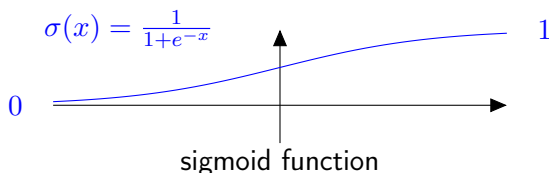


Activation function

Two linear layers is the same as one layer.

$$(\mathbf{a}^\top W_1 + \mathbf{b}_1^\top)W_2 + \mathbf{b}_2^\top = \mathbf{a}^\top W_1 W_2 + (\mathbf{b}_1^\top W_2 + \mathbf{b}_2^\top)$$

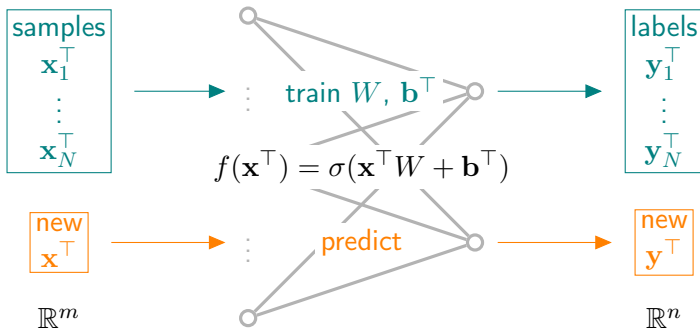
An **activation function** adds nonlinearity to the function.



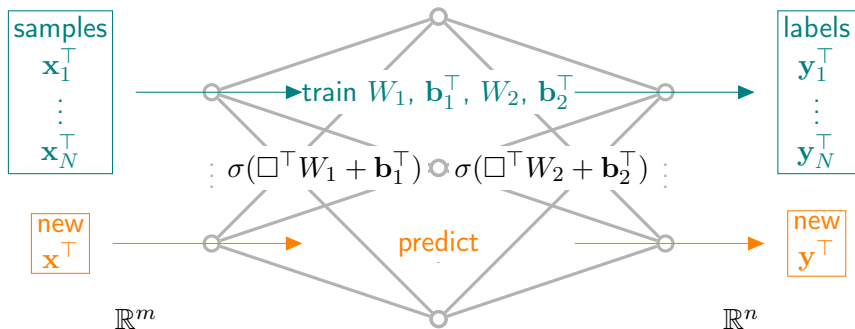
Define $\sigma(x_1, \dots, x_m) = (\sigma(x_1), \dots, \sigma(x_m))$.

So $\frac{\partial \sigma}{\partial \mathbf{x}}(\Delta_{\mathbf{x}}) = \Delta_{\mathbf{x}} \circ (\sigma'(x_1), \dots, \sigma'(x_m)) = \Delta_{\mathbf{x}} \circ \sigma'(\mathbf{x})$.

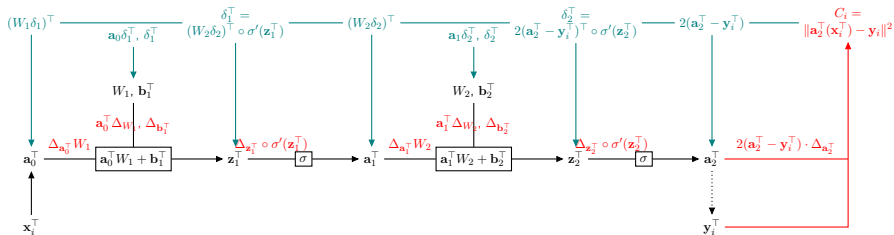
Neural network: Single-layer perceptron



Neural network: Multilayer perceptron

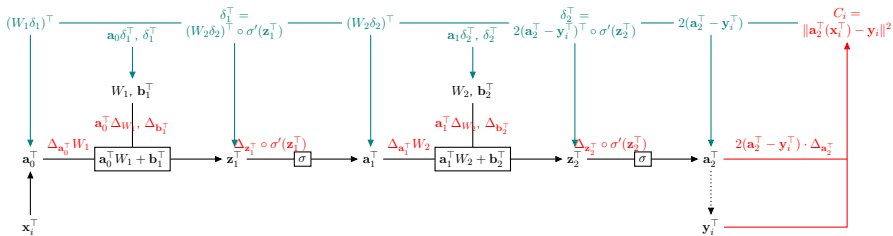


Neural network: Multilayer perceptron

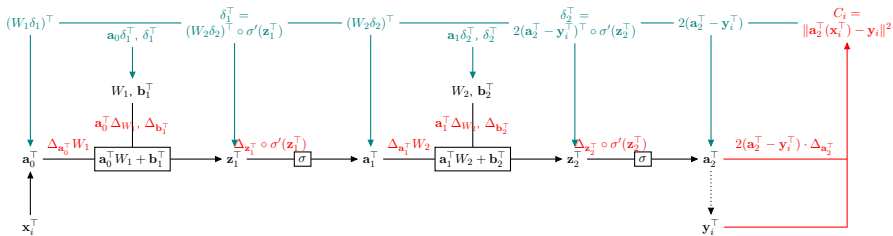


Some examples





Thanks!



Thanks!

References I



Michael A. Nielsen.

Neural Networks and Deep Learning.

Determination Press, 2015.

<http://neuralnetworksanddeeplearning.com/>.